


RESEARCH ARTICLE

Intraclass correlation: Improved modeling approaches and applications for neuroimaging

Gang Chen¹  | Paul A. Taylor¹ | Simone P. Haller² | Katharina Kircanski² |
Joel Stoddard³ | Daniel S. Pine⁴ | Ellen Leibenluft² | Melissa A. Brotman² |
Robert W. Cox¹

¹Scientific and Statistical Computing Core, National Institute of Mental Health, National Institutes of Health, Bethesda, MD

²Section on Mood Dysregulation and Neuroscience, Emotion and Development Branch, National Institute of Mental Health, Bethesda, MD

³Division of Child and Adolescent Psychiatry, Department of Psychiatry, University of Colorado School of Medicine, Aurora, Colorado

⁴Section on Development and Affective Neuroscience, Emotion and Development Branch, National Institute of Mental Health, Bethesda, MD

Correspondence

Gang Chen, Scientific and Statistical Computing Core, National Institute of Mental Health, National Institutes of Health, USA.
Email: gangchen@mail.nih.gov

Abstract

Intraclass correlation (ICC) is a reliability metric that gauges similarity when, for example, entities are measured under similar, or even the same, well-controlled conditions, which in MRI applications include runs/sessions, twins, parent/child, scanners, sites, and so on. The popular definitions and interpretations of ICC are usually framed statistically under the conventional ANOVA platform. Here, we provide a comprehensive overview of ICC analysis in its prior usage in neuroimaging, and we show that the standard ANOVA framework is often limited, rigid, and inflexible in modeling capabilities. These intrinsic limitations motivate several improvements. Specifically, we start with the conventional ICC model under the ANOVA platform, and extend it along two dimensions: first, fixing the failure in ICC estimation when negative values occur under degenerative circumstance, and second, incorporating precision information of effect estimates into the ICC model. These endeavors lead to four modeling strategies: linear mixed-effects (LME), regularized mixed-effects (RME), multilevel mixed-effects (MME), and regularized multilevel mixed-effects (RMME). Compared to ANOVA, each of these four models directly provides estimates for fixed effects and their statistical significances, in addition to the ICC estimate. These new modeling approaches can also accommodate missing data and fixed effects for confounding variables. More importantly, we show that the MME and RMME approaches offer more accurate characterization and decomposition among the variance components, leading to more robust ICC computation. Based on these theoretical considerations and model performance comparisons with a real experimental dataset, we offer the following general-purpose recommendations. First, ICC estimation through MME or RMME is preferable when precision information (i.e., weights that more accurately allocate the variances in the data) is available for the effect estimate; when precision information is unavailable, ICC estimation through LME or the RME is the preferred option. Second, even though the absolute agreement version, ICC(2,1), is presently more popular in the field, the consistency version, ICC(3,1), is a practical and informative choice for whole-brain ICC analysis that achieves a well-balanced compromise when all potential fixed effects are accounted for. Third, approaches for clear, meaningful, and useful result reporting in ICC analysis are discussed. All models, ICC formulations, and related statistical testing methods have been implemented in an open source program 3dICC, which is publicly available as part of the AFNI suite. Even though our work here focuses on the whole-brain level, the modeling strategy and recommendations can be equivalently applied to other situations such as voxel, region, and network levels.

KEYWORDS

ANOVA, estimate precision, gamma density prior for variances, intraclass correlation, linear mixed-effects modeling, reliability

1 | INTRODUCTION

Recently, reliability and reproducibility have been hot topics in science in general and in the neuroimaging community in particular. It is known that neuroimaging data are very noisy and a large proportion of the fMRI variability cannot be properly accounted for. It has been reported that less than half of the data variability can currently be explained in the typical data analysis (Gonzalez-Castillo, Chen, Nichols, & Bandettini, 2017). For example, the typical analytical approach is to make a strong and unrealistic assumption that a hemodynamic response is the same across brain regions, subjects, groups, different tasks, or conditions. In addition, even though large amounts of physiological confounding effects are embedded in the data, it remains a daunting task to fully incorporate the physiological noise in the model.

Recent surveys observed that about 60% of published experiments failed to survive replication in psychology (Baker, 2015) and about 40% in economics (Bohannon, 2016), and the situation with neuroimaging is likely to be equally, if not more, pessimistic (Griffanti et al., 2016). While in neuroimaging, reproducibility is typically a qualitative description as to whether an activation cluster is reported across many studies, the reliability of an effect estimate quantitatively describes the variation in repeated measurements performed on the same measuring entities (e.g., subjects in neuroimaging) under the identical or approximately the same experimental conditions (NIST, 2007). Specifically, reliability can be defined as the agreement or consistency across two or more measurements, and intra-class correlation (ICC) has been specifically developed for this purpose (Shrout & Fleiss, 1979; McGraw & Wong, 1996).

Generally speaking, the conventional ICC metric indicates agreement, similarity, stability, consistency, or reliability among multiple measurements of a quantity. For instance, such a quantity can be the ratings of n targets assessed by k raters or judges in a classical example (Shrout & Fleiss, 1979). In the neuroimaging context, when the same set of measuring entities (e.g., subjects) goes through the same experiment protocol under the same conditions, ICC can be utilized to assess the data quality. Those multiple measurements can be the effect estimates from n subjects under k different replications (e.g., runs, sessions, scanners, sites, twins, siblings, parent-child pairs, studies, assessments, diagnoses, or analytical methods). The same quantity (rating or effect estimate) across those k replications in the ICC definition is reflected in the word *intraclass*, as opposed to the Pearson (or *interclass*) correlation coefficient that reveals the linear relationship between two quantities that can be of different nature (e.g., brain response and cortex thickness).

Here we first review the various versions of ICC definition and their computation formulations under their classic ANOVA platform, and then we discuss their limitations and drawbacks as our motivations to develop more extended models. We then describe and validate several new, improved approaches for ICC estimation. Even though our work mainly focuses on whole-brain data analysis in neuroimaging, the methodologies can be applied to other contexts or fields when the underlying assumptions are met.

2 | VARIOUS TYPES OF ICC

In this section, we introduce the three classic types of ICC, motivating each from their basic statistical model and describing their interpretation, applicability, and generalization. Throughout this article, regular italic letters in lower case (e.g., a) stand for scalars and random variables; boldfaced italic letters in lower (\mathbf{a}) and upper (\mathbf{X}) cases for column vectors and matrices, respectively; Roman and Greek letters for fixed and random effects, respectively, on the right-hand side of a model equation. In the neuroimaging context, let y_{ij} be the effect estimate (BOLD response in percent signal change or connectivity measurement) at the i th level of within-subject (or repeated-measures) factor A and the j th level (usually subject) of factor B ($i=1, 2, \dots, k; j=1, 2, \dots, n$). When both factors A (e.g., runs, sessions, scanners, sites) and B are modeled as random effects, we have a two-way random-effects ANOVA system,

$$y_{ij} = b_0 + \pi_i + \lambda_j + \epsilon_{ij}, \quad (1)$$

where b_0 is a fixed effect or constant representing the overall average, π_i is the random effect associated with the i th level of factor A, λ_j represents the subject-specific random effect, and ϵ_{ij} is the residual. With the random variables π_i , λ_j , and ϵ_{ij} assumed to be independent and identically distributed with $N(0, \sigma_\pi^2)$, $N(0, \sigma_\lambda^2)$, and $N(0, \sigma_\epsilon^2)$, respectively, the associated ICC for the model (Equation 1) is defined as

$$ICC(2, 1) = \rho_2 = \frac{\sigma_\lambda^2}{\sigma_\pi^2 + \sigma_\lambda^2 + \sigma_\epsilon^2}, \quad (2)$$

and is numerically evaluated by

$$\hat{\rho}_2 = \frac{MS_\lambda - MS_\epsilon}{\frac{k}{n}(MS_\pi - MS_\epsilon) + MS_\lambda + (k-1)MS_\epsilon}, \quad (3)$$

where MS_π , MS_λ , and MS_ϵ are the mean squares¹ (MSs) associated with the factor A effects π_i , subject effects λ_j and the residuals ϵ_{ij} , respectively, in the ANOVA framework (Equation 1). The definition (Equation 3) is usually referred to as ICC(2,1) in the literature (Shrout & Fleiss, 1979; McGraw & Wong, 1996).

To make statistical inference, Fisher's transformation for ICC value ρ (McGraw & Wong, 1996),

$$z = \frac{1}{2} \sqrt{\frac{2(n-2)(k-1)}{k}} \ln \frac{1+(k-1)\rho}{1-\rho} \quad (4)$$

approximately follows a standard Gaussian $N(0, 1)$ under the null hypothesis $H_0: \rho = 0$, and offers a solution for significance testing. However, a better approach is to formulate an F -statistic (McGraw & Wong, 1996),

$$F_2(n-1, n(k-1)) = \frac{MS_\lambda}{MS_\epsilon}, \quad (5)$$

whose distribution is exact under $H_0: \rho_2 = 0$, unlike the Fisher transformation (Equation 4).

The meaning of the ICC can be interpreted in four common perspectives:

¹Under the ANOVA formulation, the mean squares of a factor is the sum of squares for the factor divided by the associated degrees of freedom.

- i As the definition (Equation 2) itself indicates, the ICC is the proportion of total variance that is attributed to a random factor (or accounted for by the association across the levels of the random factor). For instance, if the variance associated with subjects increases, subjects would be less similar while the levels of factor A (e.g., runs) tend to be relatively more similar, leading a higher ICC value. This proportionality interpretation is straightforwardly consistent with the non-negativity of ICC and its range of [0, 1].
- ii The ICC is the expected correlation between two effect estimates that are randomly drawn among the levels of factor A within the same level of factor B. For instance, say that ICC(2,1) for an fMRI study shows the relatedness among multiple runs. Specifically, with the assumptions in the model (Equation 1), ICC(2,1) is essentially the Pearson correlation of the effect estimates between any two levels (e.g., runs) of factor A, i_1 , and i_2 ($i_1 \neq i_2$),

$$\text{Corr}(y_{i_1j}, y_{i_2j}) = \frac{\text{Cov}(\pi_{i_1} + \lambda_j + \epsilon_{i_1j}, \pi_{i_2} + \lambda_j + \epsilon_{i_2j})}{\sqrt{\text{Var}(\pi_{i_1} + \lambda_j + \epsilon_{i_1j})\text{Var}(\pi_{i_2} + \lambda_j + \epsilon_{i_2j})}} = \frac{\sigma_\lambda^2}{\sigma_\pi^2 + \sigma_\lambda^2 + \sigma_\epsilon^2}.$$

However, it is worth emphasizing that this equivalence between ICC and Pearson correlation holds because of the following fact: ICC is a relationship gauge between, for example, any two runs in light of the same physical measure (e.g., BOLD response in ICC(2,1) as opposed to Pearson correlation between, for example, weight and height). When a run has generally a higher (or lower) effect estimate relative to the group mean effect, or when there is some extent of consistency among subjects within each run, then those effect estimates are correlated, and the ICC formulation (Equation 2) basically captures that correlation or consistency.

- i ICC is an indicator of homogeneity of the responses among the levels (e.g., subjects) of the random factor B: a higher ICC means more similar or homogeneous responses. On the other hand, when an ICC value is close to zero, the effect estimates associated with a factor A level are no more similar than those from different subjects, and the random effect components could be removed from the model (Equation 1).
- ii ICC reflects the extent of common conditions (e.g., same task and scanning parameters) that the effect estimates share. The ICC would be higher if effect estimates associated with a subject were under more similar environments.

The decision of an explanatory variable in a model as either fixed or random effects can be subtle, and the distinction is usually determined in light of the nature of the factor: interchangeability of factor levels, or whether there exists some systematic difference across the factor levels. For example, subjects are often considered as the levels of a random factor because: (a) they are each recruited through a random sampling process as representatives (or samples) of a potential or conceptual population (as embodied in the assumption of Gaussian distribution $N(0, \sigma_\pi^2)$ for the random effects λ_j in the model (Equation 1)), achieving the goal of generalization in statistical inferences; (b) their

order does not matter in the model and can be randomly permuted due to exchangeability; and (c) a particular set of subjects can, in practice, be replaced by another set. In contrast, patients and controls are typically handled as fixed effects because of the lack of exchangeability. In neuroimaging scanners or sites can be thought of as the levels of either a random- or fixed-effects factor, depending on whether the scanning parameters are similar or different across scanners or sites. Similarly, runs or sessions should be treated as random effects if no significant systematic difference exists across runs or sessions; otherwise, they should be modeled as fixed effects when habituation or familiarity effect of the task is substantial.

Because of the distinction between fixed and random effects, there is an alternative ICC definition in which the factor A (e.g., runs, sessions, scanners, sites) is modeled as fixed effects in a two-way mixed-effects ANOVA structure,

$$y_{ij} = b_0 + b_i + \lambda_j + \epsilon_{ij}, \quad (6)$$

where b_i and λ_j represent the fixed effects² of factor A and random effects of subjects (or families, in the case of parent versus child), respectively. The associated ICC is defined as

$$\text{ICC}(3, 1) = \rho_3 = \frac{\sigma_\lambda^2}{\sigma_\lambda^2 + \sigma_\epsilon^2}, \quad (7)$$

which has the same formula as Equation 2, but we note that the presence of b_i here means that the two models would have different estimates of σ_ϵ^2 . The ICC in this case can be computed as

$$\hat{\rho}_3 = \frac{MS_\lambda - MS_\epsilon}{MS_\lambda + (k-1)MS_\epsilon}, \quad (8)$$

with an exact F -statistic,

$$F_3(n-1, (n-1)(k-1)) = \frac{MS_\lambda}{MS_\epsilon}, \quad (9)$$

for significance testing under the null hypothesis $H_0: \rho_3 = 0$.

The two ICC definitions, Equations 2 and 7, are popularly notated as ICC(2,1) and ICC(3,1), respectively, and are sometimes referred to as “intertest ICCs,” extensions of the classic inter-rater reliability (Shrout & Fleiss, 1979). When there are only two levels for factor A ($k = 2$), these two versions are usually referred to as test-retest reliability measures (Zuo & Xing, 2014). When factor A represents scanning sites, these two ICC types are aligned with another term, multisite reliability, in the literature. Yet there is another ICC type in which the effects of factor A are not explicitly modeled. A prototypical example is the scenario with the effect estimates of twins from each family where there is no meaningful way to assign an order or sequence among the levels of factor A consistently among the levels of factor B (e.g., ordering twins within each family). Let y_{ij} be the effect estimate from the i th level of factor A and j th level of factor B (e.g., the i th member of family

²For simplicity, the notations for the model terms and for the corresponding variance and MS terms in the ICC formulas are undifferentiated across models. To avoid confusion, we emphasize that the values for the same notation term may change across models. For example, residuals ϵ_{ij} and its variance estimate $\hat{\sigma}_\epsilon^2$ may differ, depending on how the effects associated with factor A are modeled.

$j)$ ($i=1, 2, \dots, k; j=1, 2, \dots, n$). A one-way random-effects ANOVA can be formulated to decompose the response variable or effect estimate y_{ij} as

$$y_{ij} = b_0 + \lambda_j + \epsilon_{ij}, \quad (10)$$

where λ_j codes the random effect of the j th level of factor B (e.g., family j).

The ICC for the model in Equation 10 is defined as

$$ICC(1, 1) = \rho_1 = \frac{\sigma_\lambda^2}{\sigma_\lambda^2 + \sigma_\epsilon^2}, \quad (11)$$

which can be estimated similarly as (Equation 8),

$$\hat{\rho}_1 = \frac{MS_\lambda - MS_\epsilon}{MS_\lambda + (k-1)MS_\epsilon}, \quad (12)$$

where MS_λ and MS_ϵ are the mean squares (MSs) associated with the family effects λ_j and the residuals, respectively, in the ANOVA framework of model (Equation 10). The testing statistic for Equation 11 has the form as Equation 5, and the definition (Equation 11) is usually referred to as ICC(1,1) in the literature (Shrout & Fleiss, 1979; McGraw & Wong, 1996).

The four interpretation perspectives for ICC(2,1) also apply directly to other two types. For example, with the assumptions in the model (Equation 6), ICC(3,1) is a special case of Pearson correlation of the effect estimates between any two levels of factor A, i_1 , and i_2 ($i_1 \neq i_2$),

$$\text{Corr}(y_{i_1j}, y_{i_2j}) = \frac{\text{Cov}(\lambda_j + \epsilon_{i_1j}, \lambda_j + \epsilon_{i_2j})}{\sqrt{\text{Var}(\lambda_j + \epsilon_{i_1j})\text{Var}(\lambda_j + \epsilon_{i_2j})}} = \frac{\sigma_\lambda^2}{\sigma_\lambda^2 + \sigma_\epsilon^2}.$$

Nevertheless, ICC(3,1) is more similar to Pearson correlation than ICC(1,1) and ICC(2,1) in the sense that each level of factor A is assumed to have a different mean, but it remains unlike Pearson correlation as the assumption of same variance holds across the levels of factor A for all the three ICC definitions. In addition, ICC(2,1) and ICC(3,1) are sometimes called *absolute* agreement and *consistent* agreement, respectively. The distinction between these two ICC types can be hypothetically illustrated by paired effect estimates (e.g., in percent signal change for an fMRI experiment) from five subjects during two sessions (Table 1). The three ICC values are all nonnegative because they represent a proportion of total variance embedded in the data, and they generally follow a sequential order (Shrout & Fleiss, 1979): $ICC(1, 1) \leq ICC(2, 1) \leq ICC(3, 1)$.

The first index in the ICC(\cdot, \cdot) notation specifies the ICC type, while the second indicates the relationship between two *single* measurements (e.g., between two twins for ICC(1,1) or two levels of factor A for ICC(2,1) and ICC(3,1)). For each of the three *single measurement* ICC types, there is another version, called *average measurement* ICC, which shows the relationship between two sets of average measurements among the k levels of factor A, and with notations ICC(1, k), ICC(2, k), and ICC(3, k), they are similarly defined as the single measurements version except that the terms in the denominator, σ_ϵ^2 , $\sigma_\pi^2 + \sigma_\epsilon^2$, and σ_π^2 , for ICC(1,1), ICC(2,1), and ICC(3,1), respectively, are each scaled by a factor of k^{-1} . By definition, the average measurement ICC is larger than its single measurements counterpart. In addition, a similar correlation interpretation about the average measurement ICCs can be seen with, for example, ICC(2, k),

TABLE 1 Hypothetical effect estimates (e.g., BOLD response in percent signal change) with a relationship from five subjects during two sessions ($y_{2j} = y_{1j} + 0.2, j=1, 2, \dots, 5$; $ICC(1,1) = 0.43$; $ICC(2,1) = 0.56$; $ICC(3,1) = 1$, Pearson correlation $r = 1$)

Subject	s1	s2	s3	s4	s5
Session 1, y_{1j}	0.1	0.2	0.3	0.4	0.5
Session 2, y_{2j}	0.3	0.4	0.5	0.6	0.7

$$\begin{aligned} \text{Corr}\left(\frac{1}{k} \sum_{i=1}^k y_{ij}^{(1)}, \frac{1}{k} \sum_{i=1}^k y_{ij}^{(2)}\right) &= \frac{\text{Cov}\left(\frac{1}{k} \sum_{i=1}^k \pi_i^{(1)} + \lambda_j + \frac{1}{k} \sum_{i=1}^k \epsilon_{ij}^{(1)}, \frac{1}{k} \sum_{i=1}^k \pi_i^{(2)} + \lambda_j + \frac{1}{k} \sum_{i=1}^k \epsilon_{ij}^{(2)}\right)}{\sqrt{\text{Var}\left(\frac{1}{k} \sum_{i=1}^k \pi_i^{(1)} + \lambda_j + \frac{1}{k} \sum_{i=1}^k \epsilon_{ij}^{(1)}\right) \text{Var}\left(\frac{1}{k} \sum_{i=1}^k \pi_i^{(2)} + \lambda_j + \frac{1}{k} \sum_{i=1}^k \epsilon_{ij}^{(2)}\right)}} \\ &= \frac{\sigma_\lambda^2}{\frac{1}{k} \sigma_\pi^2 + \sigma_\lambda^2 + \frac{1}{k} \sigma_\epsilon^2}, \end{aligned}$$

where the superscript such as those in $y_{ij}^{(1)}$ and $y_{ij}^{(2)}$ indicates a particular set of data substantiation, and thus $\pi_i^{(m)} \sim \text{iid } N(0, \sigma_\pi^2)$, $\lambda_j^{(m)} \sim \text{iid } N(0, \sigma_\lambda^2)$, $\epsilon_{ij}^{(m)} \sim \text{iid } N(0, \sigma_\epsilon^2)$. As the average measurement ICC is less popular in practice, we hereafter focus on their single measurement counterpart, but our modeling work below can be directly extended to the average measurements ICC.

3 | LITERATURE SURVEY OF ICC FOR NEUROIMAGING

ICC has been applied to neuroimaging data for over 10 years, mainly to examine reliability under various scenarios. In particular, ICC(2,1) has been largely adopted in the field, using the ANOVA approach. ICC(2,1) has been used to show reliability under various scenarios, for example, at the regional level (Fiecas et al., 2013), at the network level for resting-state (Cao et al., 2014; Guo et al., 2012), at the whole-brain level with ANOVA (Kristo et al., 2014; Quiton, Keaser, Zhuo, Gullapalli, & Greenspan, 2014; Zanto, Pa, & Gazzaley, 2014) and with linear mixed-effects (LME) modeling using a precursor of 3dICC in AFNI (Fiecas et al., 2013; Haller et al., 2017; White et al., 2016). ICC(3,1) has been applied at the regional level for task-related fMRI data (Cáceres, Hall, Zelaya, Williams, & Mehta, 2009; Jaeger et al., 2015) and for MEG data (Recasens & Uhlhaas, 2017), at the network level for resting-state data (Braun et al., 2012), on the regional homogeneity of resting-state data (Zuo et al., 2013), and at the whole-brain level for task-related fMRI data, using ANOVA in PASW and/or Matlab (Brandt et al., 2013; Cáceres et al., 2009) and in SAS (Fournier, Chase, Almeida, & Phillips, 2014). It should be noted that in many cases the ICC type adopted in a study was not clearly stated, often due to the ambiguity in terms of the model involved (Lin et al., 2015; Shah, Cramer, Ferguson, Birn, & Anderson, 2016; Zuo et al., 2010b).

There have been occasions in which ICC(1,1) and ICC(3,1) have been explicitly employed at times in the literature. For example, ICC(1,1) has been applied to brain networks based on resting-state fMRI data (Wang et al., 2011), to functional near-infrared spectroscopy

TABLE 2 A hypothetical study of 17 subjects with missing data (1: available; 0: missing)

Subject	s1	s2	s3	s4	s5	s6	s7	s8	s9	s10	s11	s12	s13	s14	s15	s16	s17
Session 1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
Session 2	1	1	1	1	0	1	1	0	1	1	1	1	1	1	1	1	1
Session 3	1	0	1	1	1	0	1	1	1	1	1	0	1	1	0	1	1

(fNIRS) data (e.g., Bhambhani, Maikala, Farag, & Rowland, 2006; Plichta et al., 2006, 2007; Tian et al., 2012; Zhang et al., 2011), and to resting-state data at the whole-brain level (Zuo et al., 2010a), and to task-related data through LME at the regional level (Töger et al., 2017). However, we note that in a large number of publications, the adoption of the ICC type was neither explicitly explained nor justified, which makes precise interpretation difficult.

On the whole-brain level, in addition to the ICC computation through LME as implemented in the open source AFNI program 3dLME (Chen, Saad, Britton, Pine, & Cox, 2013) and in DPARBI (Yan, Wang, Zuo, & Zang, 2016), there have been a few Matlab toolboxes publicly available: three using ANOVA (Cáceres et al., 2009; Fiecas et al., 2013; Molloy & Birn, 2014), and one using both ANOVA and LME in Matlab (Zuo et al., 2010b). The concept of ICC has also been extended to characterize cross-subject heterogeneity when the effect estimate precision is incorporated into fMRI group analysis (Chen, Saad, Nath, Beauchamp, & Cox, 2012), and to reveal the relatedness among subjects in intersubject correlation (ISC) analysis with data from naturalistic scanning (Chen, Taylor, Shin, Reynolds, & Cox, 2017a).

To address the reliability and reproducibility issue on a large scale, the Consortium for Reliability and Reproducibility (CoRR) has been established to explore the test-retest reliability as a minimum standard for methods development in functional connectomics (Zuo et al., 2014). With thousands of datasets openly available, reliability could be tested through a variety of perspectives. One specific such effort was demonstrated in estimating both intra- and cross-subject variability as well as the reliability of multiple metrics with a resting-state study of 10 sessions across 30 subjects (Chen et al., 2015). A multisite resting-state study (Noble et al., 2017) recently showed that a relatively poor reliability with a short scan (e.g., 5 min) could be improved to some extent with a longer scan (e.g., 25 min). A low reliability for resting-state data was also seen with an extended version of ICC proposed as a global measure of reliability across an ensemble of ROIs (Shou et al., 2013).

4 | MOTIVATIONS FOR FURTHER MODELING WORK

The ANOVA framework has been adopted in computing ICC through the MS terms largely for historical reasons because ANOVA was developed and became widely adopted in early twentieth century: the framework is widely introduced in basic statistics textbooks, and the MS terms are efficient to compute. Various computational tools are widely available through ANOVA (e.g., packages *irr* and *psych* in R).

However, the ANOVA approach does have both limitations for interpretation and practical drawbacks for calculation:

- i *Negative value.* Although the ICC value should be theoretically nonnegative per its definition as the proportionality of the total variance (as in Equation 2), its estimation (as in Equation 3), may become negative due to the fact that the numerator in the computational formula is the difference between two MS terms. Such cases are uninterpretable. Importantly, in neuroimaging negative ICCs are not rare occurrences, with a large number of such voxels appearing in the brain (and in any tissue type).
- ii *Missing data.* In common practical situations, missing data may occur. As data balance is essential in partitioning the MS terms, ANOVA cannot properly handle missing data due to the breakdown of the underlying rigid variance-covariance structure. For example, the data for the six subjects who missed scanning for one session as shown in Table 2 would have to be abandoned due to the rigidity of the ANOVA structure.
- iii *Confounding effects.* Under some circumstances it might be desirable to incorporate explanatory variables into a model, so that the variability due to those confounding effects can be properly accounted for. For example, subject groupings (e.g., sex, handedness, genotypes), age, and reaction time are typical candidates in neuroimaging that often warrant their inclusion in the model. However, the rigid ANOVA structure usually does not easily allow for such inclusions.³
- iv *Sampling errors.* Conceptually, the residual term ϵ_{ij} in the ICC model can be considered to represent measurement or sampling errors. In other words, the underlying assumption is that all the effect estimates y_{ij} share the same sampling variance for the measurement errors. However, unlike the typical situation in other fields where the effect estimates are usually direct measurements (e.g., scores assessed by raters), the effect estimates in neuroimaging are generally obtained through a data reduction process with a regression model for each subject. That is, each effect estimate is associated with a standard error that indicates the precision⁴ of the estimation, and heterogeneity or (possibly unfounded) heteroscedasticity is expected because the standard error varies across subjects and between twins or parents/

³One exception is that, when a subject-grouping factor (e.g., males versus females) is considered, it is possible to construct the involved MS terms in the special case of having exactly equal number of subjects across all groups. However, even for such a balanced scenario, specific MS terms would have to be derived for each ICC computation formula.

⁴Precision is defined as the reciprocal of the variance.

children, across runs/sessions, or scanners/sites. When the standard error for the effect estimate is ignored in favor of a homoscedasticity assumption in the ANOVA formulation (as widely practiced in neuroimaging when computing the ICCs, and in group analysis), it raises the question: what is the impact for the ICC estimate when heterogeneity of standard error is substantially present across the measuring entities?

- v *Type selection.* The applicability of ICC(1,1) is limited to situations where the levels of the repeated-measures factor are measuring entities that are difficult to assign meaningful orders or sequences such as twins. However, between ICC(2,1) and ICC(3,1), the choice becomes challenging for the investigator, with nontrivial impact on the ICC results: other than some prior knowledge about the potential existence of confounding effects or systematic difference across the levels of the repeated-measures factor, there is no solid statistical tool under ANOVA to leverage one choice over the other. For example, is there any statistical evidence that could allow us to decide unequivocally between ICC(2,1) and ICC(3,1) by treating runs or sessions as random or fixed effects? In addition, the typical whole-brain analysis through a massively univariate approach at the voxel level may further aggravate the choice.

Our ICC modeling work here hinges around these five limitations of the ANOVA approach. We first discuss four alternative modeling approaches as extensions to the ANOVA framework, and then use an experimental dataset to examine the performance of the various modeling methods. Each of the four models addresses the limitations and drawbacks we discussed above, and provides incremental improvements (Table 3). Further discussion is presented at the end. The implementations of our modeling work are publicly available for voxel-wise computation through program 3dICC as part of the AFNI suite (Cox, 1996). As ICC(1,1) is typically adopted in neuroimaging for special cases of studies with twins, our focus here is on the other two types due to their wider applicability. Nevertheless, the modeling strategies discussed here can be directly expanded to ICC(1,1), and possibly can be further applied to other fields, when appropriate, even though our focus remains on neuroimaging.

5 | THEORY: THREE EXTENDED ICC MODELS

Here we propose four mixed-effects models: linear mixed-effects (LME), regularized mixed-effects (RME), multilevel mixed-effects (MME), and regularized multilevel mixed-effects (RMME). These four models are introduced in a sequential order, reflecting their incremental improvements.

5.1 | Linear mixed-effects (LME) modeling

Whenever multiple values (e.g., effect estimates from each of two scanning sessions) from each measuring entity (e.g., subject or family) are correlated (e.g., the levels of a within-subject or repeated-measures factor), the data can be formulated with an LME model, sometimes also

TABLE 3 Issues in ICC computations under ANOVA, and summary of which mixed-effects models can (✓) or cannot (✗) address each

Issues	ANOVA	LME	RME	MME	RMME
Negative ICC	✗	✓	✓	✓	✓
Zero ICC	✗	✗	✓	✗	✓
Missing data	✗	✓	✓	✓	✓
Confounding effects	✗	✓	✓	✓	✓
Sampling error	✗	✗	✗	✓	✓
Type selection	✗	✓	✓	✓	✓

referred to as a multilevel or hierarchical model. One natural extension to the ANOVA modeling in Equations 1, 6, and 10 is to simply treat the model conceptually as LME, reformulating neither the equations nor their ICC definitions. This LME approach for ICC has previously been implemented in the program 3dLME (Chen et al., 2013) for voxel-wise data analysis in neuroimaging. For example, the ICC(2,1) model is an LME case with two crossed random-effects terms, whose applications can be seen under other circumstances such as intersubject correlation analysis (Chen et al., 2017a) and psycholinguistic studies (e.g., Baayen et al., 2008).

However, the application of LME methodology to ICC does not stop at the conceptual level, and in fact it has several advantages in some aspects of computation where limitations are present under the ANOVA framework. Specifically, the variances for the random effects components and the residuals are directly estimated through optimizing the restricted maximum likelihood (REML) function, and thus the ICC value is computed with variance estimates $\hat{\sigma}_\pi^2$, $\hat{\sigma}_\lambda^2$, and $\hat{\sigma}_\epsilon^2$, through the definitions in Equations 11, 2, 7, instead of with their counterparts with MS terms, (Equations 12, 3, 8), under ANOVA. Therefore, in conjunction with the theoretical quantities, the estimated ICCs are non-negative by definition, avoiding the interpretability difficulties that ANOVA-based estimates can present when negative. Similarly, the two *F*-statistic formulas (Equations 5 and 9) can be expressed in terms of variance estimates as well,

$$F = \frac{k\hat{\sigma}_\lambda^2}{\hat{\sigma}_\epsilon^2} + 1. \quad (13)$$

Another convenient byproduct from the LME model interpretation of (Equation 1) for ICC(2,1) is that the ICC for the absolute agreement between any two measuring entities (e.g., subjects) can be easily obtained through ratio among the variances,

$$\hat{\rho}_2 = \frac{\hat{\sigma}_\pi^2}{\hat{\sigma}_\pi^2 + \hat{\sigma}_\lambda^2 + \hat{\sigma}_\epsilon^2},$$

which is usually not discussed under the ANOVA framework.

In regard to the type selection choice between ICC(2,1) and ICC(3,1), which is an ambiguous one to some extent in the ANOVA framework, in the LME model (Equation 6), the fixed effects associated with the levels of the repeated-measures factor *A* can now be directly examined through statistical assessment: If the fixed effects b_i can be deemed negligible (i.e., no confounding effects nor systematic

differences across the factor A levels), then ICC(2,1) offers a better metric because the model (Equation 1) is more parsimonious. We will elaborate this point later through an experimental dataset.

Furthermore, missing data can be naturally handled in LME because parameters are estimated through the optimization of the (restricted) maximum likelihood function, where a balanced structure is not required. As long as the missing data can be considered to occur randomly without structure (i.e., no systematic pattern exists among the missing data), and there are enough measuring entities present (e.g., $n \geq 10$) for each level of the repeated-measures factor, then the computation can still be performed. Specifically, if no relationship exists between whether a data point is missing and any values in the data set (missing or observed), the situation is considered missing completely at random (MCAR); missing at random (MAR) occurs if the missingness is not fully random (e.g., men being more likely to participate in fMRI scanning), but it can be fully controlled by relevant variables (e.g., incorporating sex as a covariate). When a situation with MCAR or MAR occurs, the conventional ANOVA or GLM cannot handle missing data because of the loss of the rigid variance-covariance structure while LME does not rely on such a rigid structure.

In addition, the extension to incorporate confounding effects is readily available through adding more fixed-effects terms into the model. For instance, the ICC(2,1) model (Equation 1) can be expanded to an LME model with two crossed random-effects components,

$$y_{ij} = b_0 + \sum_{l=1}^m b_l x_j^{(l)} + \pi_i + \lambda_j + \epsilon_{ij}, \quad (14)$$

where $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ are m explanatory variables (e.g., sex, age) that can be either categorical or quantitative, and b_1, b_2, \dots, b_m are their corresponding fixed effects.

For the convenience of further discussion, we adopt the conventional LME platform for ICC with a more generic and inclusive formulation (Pinheiro and Bates, 2004; Chen et al., 2017a) than (Equations 10, 1, 6). The following formulation contains all of the aforementioned models as special cases, and will be further discussed hereafter:

$$y = Xb + Z\theta + \epsilon, \quad (15)$$

where the known vector⁵ $y_{kn \times 1} = \text{vec}(y_{ij})$ is the vectorization or column-stacking of effect estimates y_{ij} ; $b_{(m+1) \times 1}$ contains the unknown fixed-effects to be estimated; the known model (or design) matrix $X_{kn \times (m+1)}$ codes the explanatory variables for the fixed effects; the known model (or design) matrix Z is usually a column-wise subset of X ; θ contains the unknown random effects to be estimated; and $\epsilon_{kn \times 1}$ contains the unknown residuals. The distributional assumptions are that $\theta \sim N(0, V)$, and $\epsilon \sim N(0, R)$, where V and R are the unknown variance-covariance matrices for the random effects θ and residuals ϵ , respectively; also, θ and ϵ are independent of each other (i.e., $\text{Cov}(\theta, \epsilon) = 0$). For the ICC context with no missing data, $R = \sigma_\epsilon^2 I_{kn}$; for ICC(1,1) and ICC(3,1), $V = \sigma_\lambda^2 I_{kn}$, and for ICC(2,1), $V = I_n \otimes \text{diag}(\sigma_\pi^2, \sigma_\lambda^2)$.

For those ANOVA cases where a negative ICC value would occur due to the subtraction between two MS terms, LME avoids negativity by having a lower boundary at 0, via a positive definiteness for the variance-covariance matrix when estimating the variance components through optimizing the likelihood function within the nonnegative domain of the variance components. With these considerations in mind, the question for those ambiguous values in the LME framework becomes: are those effects within a region fully uncorrelated across the levels of the repeated-measures factor; or is the zero variance estimate simply some artifact providing a zero; or is it a consequence of convergence failure in the optimization algorithms when solving LME? This is addressed by introducing an improvement to the conventional LME by regularizing the variance components. We note that all of the advantages of LME over ANOVA discussed above also carry over to the other variants of mixed-effects models below.

5.2 | Regularized mixed effects (RME) model

When a zero variance estimate occurs, typically the corresponding likelihood profile is decreasing or close to a flat line at zero. Such a scenario may occur when the sample size is small or when the signal-to-noise ratio is low, derailing the LME capability to provide a meaningful variance estimate in this near-zero boundary value of zero. It is unfortunately not rare to have either a small number of subjects or a neuroimaging signal submerged with noise. Compared to a negative ICC value (or negative variance estimate σ_λ^2), having a floor for σ_λ^2 at 0 avoids an uninterpretable situation; however, even the zero estimate for ICC may be questionable to some extent: do we truly believe that all the subjects in a study have *exactly* the same BOLD response or average effect across the levels of factor A, as implied by such a value?

Variances in LME are estimated by optimizing REML, which is equivalent to the posterior marginal mode, averaging over a uniform prior on the fixed effects. To pull out of the trapping area surrounding the likelihood boundary, one possibility is to apply a weakly informative prior distribution for the variance parameters. This regularization approach can be conceptualized as forming a compromise between two competing forces: the anchoring influence of the prior information and the strength of the data. With a reasonable prior distribution, one may prevent a numerical boundary estimate from occurring by “nudging” the initial variance estimate by no more than one standard error, leading to negligible influence from the prior when the information directly from the data should be fully relied upon (Chung, Rabe-Hesketh, Dorie, Gelman, & Liu, 2013).

Here we adopt a weakly informative prior distribution, gamma density function (Chung et al., 2013), for a variance component v in the LME model (i.e., σ_π^2 and σ_λ^2 in Equation 1, σ_λ^2 in Equation 6),

$$h(v; \eta, \kappa) = \frac{\kappa^\eta}{\Gamma(\eta)} v^{\eta-1} e^{-\kappa v}, \quad \eta, \kappa > 0, \quad (16)$$

where η and κ are the shape and rate (inverse scale) parameters, respectively. In practice, the gamma density function has two desirable properties for our nudging purpose here: (a) a positive and constant derivative at zero when $\eta = 2$, guaranteeing a nudge toward the positive direction, and (b) $h(0; \eta, \kappa) = 0$ when $\eta > 1$, allowing the variance to

⁵Without loss of generality, the dimensions shown here for vectors and matrices are assumed to have no missing data, by default. When missing data occur, the dimensions can be adjusted accordingly.

hit the boundary value of zero when the true value is zero. With $\eta = 2$, the prior is uninformed (and improper) within $(0, \infty)$ when κ approaches 0, and becomes gradually informative (but still weak) when κ is away from 0 (Chung et al., 2013). Therefore, a parameter set of $\eta = 2$ with a small rate parameter κ produces a positive estimate for the variance but does not overrule the data itself.

Typical Bayesian methodology involves estimating the posterior distribution through sampling with simulations (e.g., Markov chain Monte Carlo). However, here the prior density can be directly incorporated into the likelihood function for LME because of the conjugacy of the exponential families, hence no simulations are required to nudge the variance estimate out of the boundary value for the ICC computation. As an added benefit, there is usually little extra computational cost added to the classic LME computations. In fact, an overall higher efficiency can be achieved in some circumstances, because the prior may speed up the convergence process otherwise stuck or slowed in the trapping area close to the boundary.

With both classic LME and its “nudging” version, RME, there remains one last limitation mentioned in the Introduction to be overcome for ICC estimation: How can the LME model utilize the precision information (i.e., standard error) associated with the effect estimates from the individual subject analysis? Would the precision information provide a more accurate partitioning among the variance components including situations when ANOVA renders negative ICC or when LME forces σ_ϵ^2 to be 0?

The effect estimates of interest in neuroimaging are usually not raw data or direct measures, but instead regression coefficients as output from a time series regression model; as such, these effect estimates have their own measurement uncertainties or confidence intervals. A close inspection of the LME model for ICC (Equations 10, 1, 6) reveals that the residual term ϵ_{ij} represents the uncertainty embedded in the effect estimates. An underlying assumption in the classic LME model holds that the measurement errors share the same uncertainty: $\epsilon_{ij} \sim \text{iid } N(0, \sigma_\epsilon^2)$. As the effect estimates come from a time series regression model separately for each subject, their precision is not necessarily the same and may vary significantly across subjects for various reasons, including variations in trial sequence, data points, and outliers. Therefore, in practice, each effect estimate is associated with a different variance for its measurement error: $\epsilon_{ij} \sim N(0, \sigma_{ij}^2)$. The true value of σ_{ij}^2 is usually unknown, but its estimation is readily available, conveniently embedded in the denominator of the t-statistic for each effect estimate out of the individual subject analysis. Such an approach has previously been developed for simple group analysis, with the standard error for the effect estimate incorporated into the group model in neuroimaging (e.g., FLAME in FSL (Woolrich, Behrens, Beckmann, Jenkinson, & Smith, 2004); 3dMEMA (Chen et al., 2012)). Here we apply the same approach to ICC computation for the situation when precision information is available.

5.3 | Multilevel mixed-effects (MME) modeling

Here we extend the LME model (Equations 10, 1, 6), through replacing the assumption $\epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$ with an unknown parameter σ_ϵ^2 in the

model by $\epsilon_{ij} \sim N(0, \sigma_{ij}^2)$, where the variance estimate $\hat{\sigma}_{ij}^2$ for σ_{ij}^2 is assumed to be available, which is usually the case for task response. We call this approach multilevel mixed-effects (MME) modeling, with the term *multilevel* reflecting the fact that the modeling approach borrows part of a methodology typically adopted in robust meta-analysis when summarizing across previous studies, each of which provided both effect estimate and its standard error. The MME counterpart for the standard LME formulation (Equations 15) is extended to have the assumption $\mathbf{R} = \text{diag}(\text{vec}(\hat{\sigma}_{ij}^2))$.

Although computationally more sophisticated due to the involvement of more than one variance component in the case of the model (Equation 1) for ICC(2,1), the basic numerical scheme remains similar to our previous work for group analysis (Chen et al., 2012). That is, the variance components for the random effects are iteratively solved through optimizing REML, with the estimates $\hat{\sigma}_{ij}^2$ for measurement precision playing a weighting role: an effect estimate y_{ij} with a smaller (or larger) $\hat{\sigma}_{ij}^2$ has a larger (or smaller) impact on the estimation of the variance components and fixed effects. It is this strategy of differential weighting that separates MME from the previous models in which each of the effect estimates is treated equally.

One adjustment for MME specific to the ICC case is the following. The variance for the residuals, σ_ϵ^2 , in the ICC definitions under all other models is no longer available, due to the replacement of the residual term with an unknown variance by the measurement error with an estimated variance. In its place, we substitute $\hat{\sigma}_\epsilon^2$ with the weighted (or “typical”) average $(\hat{\sigma}_\epsilon^2)_W$ —instead of the arithmetic average—of the sample variances $\hat{\sigma}_{ij}^2$ in light of the generic model (Equation 15) (Higgins & Thompson, 2002; Viechtbauer, 2010), where

$$(\hat{\sigma}_\epsilon^2)_W = \frac{T - p}{\text{tr}(\mathbf{W} - \mathbf{W}\mathbf{X}(\mathbf{X}^T\mathbf{W}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{W})}, \quad (17)$$

and T is the total number of data points in \mathbf{y} (where $T = kn$ if no missing data occur); p is the column rank of \mathbf{X} ; $\text{tr}()$ denotes the trace operator; \mathbf{X} and $\mathbf{W} = \mathbf{R}^{-1} = \text{diag}(\frac{1}{\hat{\sigma}_{11}^2}, \dots, \frac{1}{\hat{\sigma}_{kn}^2})$ are the model matrix for the fixed effects in the model (Equation 15) and the weighting matrix, respectively. The differential weighting is reflected in the heterogeneous diagonals of the variance-covariance matrix \mathbf{R} for the measurement errors, which reduces (Equation 17) to a simplified form (Higgins & Thompson, 2002) of $(\hat{\sigma}_\epsilon^2)_W$ in the case of no missing data and no explanatory variables, that is, $\mathbf{X} = \mathbf{1}_{kn \times 1}$, in the model (Equation 15),

$$(\hat{\sigma}_\epsilon^2)_W = \frac{kn - 1}{\sum_{ij} w_{ij} - \frac{\sum_{ij} w_{ij}^2}{\sum_{ij} w_{ij}}} = \frac{1}{\frac{1}{kn-1} \left(\Omega - \sum_{ij} \frac{w_{ij}}{\Omega} w_{ij} \right)}, \quad (18)$$

where the weights $w_{ij} = \frac{1}{\hat{\sigma}_{ij}^2}$, and the total precision $\Omega = \sum_{ij} w_{ij}$, $i = 1, 2, \dots, k$, $j = 1, 2, \dots, n$. The expression for $(\hat{\sigma}_\epsilon^2)_W$ in Equation 18 can be intuitively interpreted as the following: the denominator is the difference between the total precision and the weighted mean precision, scaled by the degrees of freedom, $kn - 1$.

5.4 | Regularized MME (RMME)

A zero variance estimate may still occur under MME for the same reason as in LME, namely that the corresponding likelihood profile peaks

at zero. Similarly, the same gamma density function (Equation 16) can be adopted as a weakly informative prior distribution (Chung et al., 2013), for a variance component v in the MME model (e.g., see the example in Equation 16 here), parallel to the extension from LME to RME. Here, as well, since no posterior samplings are involved in the process, the prior sometimes can even speed up the convergence process compared to MME. We note that RMME solves all the issues raised here, with the current ANOVA framework, as summarized in Table 3.

6 | PERFORMANCE COMPARISONS AMONG THE MODELS

6.1 | Model implementations

Here, we have discussed three types of ICC, each of which can be estimated through five modeling strategies (ANOVA, LME, RME, MME, and RMME), leading to a total of 15 scenarios with an accompanying F -statistic defined by either Equation 5. These 15 models and the corresponding F -statistics are all implemented in a publicly available program 3dICC in AFNI, making use of the *R* packages *psych* (Revelle, 2016), *lme4* (Bates, Maechler, Bolker, & Walker, 2015), *blme* (Chung et al., 2013), and *metafor* (Viechtbauer, 2010). The program 3dICC takes a three-dimensional dataset as input for whole-brain voxel-wise analysis, but it can also be utilized to analyze a single voxel, region, or network. Additionally, parallel computing is available with multiple CPUs through the *R* package *snow* (Tierney, Rossini, Li, & Sevcikova, 2016). For each ICC model except ANOVA, the fixed effects (intercept or group average effect for each of the three models, as well as additional comparisons among factor A levels for ICC(3,1)) and their corresponding t -statistics are provided in the output from 3dICC, in addition to ICC and the associated F -statistic value.⁶

6.2 | Experimental testing dataset

To demonstrate the performances of our four proposed modeling approaches in comparison to ANOVA, we utilize experimental data from a previous fMRI study (Haller et al., 2017). Briefly, 25 healthy volunteers (mean age = 13.97 years, SD = 2.22 years, range = 10.04–17.51 years; 60% female) were asked to judge the gender of happy, fearful, and angry face emotions. Each emotion was displayed at three intensities (50%, 100%, and 150%). A neutral condition, representing 0% intensity, was included for each face emotion (i.e., three neutral subsets were created, one for each face emotion). MRI images were acquired in a General Electric 3T scanner (Waukesha, WI, USA), and the participants completed two MRI scanning sessions approximately two-and-a-half months apart (mean = 75.12 days, SD = 15.12 days, range: 47–109 days). The fMRI echoplanar images (EPIs) were collected with the following scan parameters: flip angle = 50°, echo time = 25

ms, repetition time = 2300 s, 47 slices, planar field of view = 240 × 240 mm², acquisition voxel size = 2.5 × 2.5 × 3 mm³, and three runs with 182 volumes for each in a total acquisition time of 21 min. The parameters for the anatomical MPRAGE images were: flip angle = 7°, inversion time = 425 ms, and acquisition voxel size = 1 mm isotropic.

The EPI time series went through the following preprocessing steps in AFNI: de-spiking, slice timing and head motion corrections, affine alignment with anatomy, nonlinear alignment to a Talairach template TT_N27, spatial smoothing with a 5 mm full-width half-maximum kernel and scaling by the voxel-wise mean. Individual TRs and the immediately preceding volume were censored if (a) the motion shift (defined as Euclidean norm of the derivative of the translation and rotation parameters) exceeded 1 mm between TRs; or (b) more than 10% of voxels were outliers.⁷ Only trials with accurate gender identification were included in the final analysis, but incorrect trials were also modeled as effects of no interest. Separate regressors were created for each of 13 event types (i.e., 30 trials for each face emotion at each intensity, neutral trials represented by three regressors of 30 trials each, and incorrect trials). Six head motion parameters and baseline drift using third order Legendre polynomials were included as additional regressors. The two sessions were analyzed separately, but the three runs with each session were concatenated⁸ and then entered into a time series regression model with a serial correlation model of ARMA(1,1) for the residuals.

The effect estimate in percent signal change, combined with the variance for the corresponding measurement errors, for the neutral condition associated with angry face-emotion from the individual subject analysis, was adopted for comparing the five ICC models: ANOVA, LME, RME, MME, and RMME. ICC(1,1) is not applicable in this case, but both ICC(2,1) and ICC(3,1) and their F -statistics were computed for each of the five models; the session effect and the corresponding t -statistic were also examined. The runtime was about one hour for each of the analyses with 16 CPUs on a Linux system (Fedora 14) with Intel® Xeon® X5650 at 2.67 GHz.

6.3 | Model comparisons

ANOVA renders a substantial number of voxels with negative ICC values (first column in Panel A, Figure 1; first row and first column in Figure 3), while LME provides virtually the same ICC estimates as ANOVA, with primary difference that those negative ICC estimates are replaced with 0 (uncolored voxels in the second column in Panel B, Figure 1; scatterplot cells (1, 2) and (2, 1) in Figure 3). It is worth noting that a significant proportion of voxels with negative or zero ICC from ANOVA or LME appear in gray matter. For RME, we tested four different priors by varying the rate parameter κ at values of 0, 0.1, 0.3, and

⁶The F -statistic is not exact in the cases of RME, MME, and RMME. However, it is important to note that the F -statistic would be an approximation too even for ICC(1,1) and ICC(3,1) under ANOVA and LME since the measurement errors are ignored.

⁷An outlier is defined as a data point where the distance between its value and the overall mean (after trend removal) is above a threshold $\sqrt{\frac{2}{T}}q^{-1}(\frac{0.001}{T})m$, where $q(x)=1-p(x)$ is the reversed Gaussian density function, T is the length of time series, and m is the median absolute deviation of the time series.

⁸Despite the concatenation, the discontinuities across runs were properly handled (Chen et al., 2012).

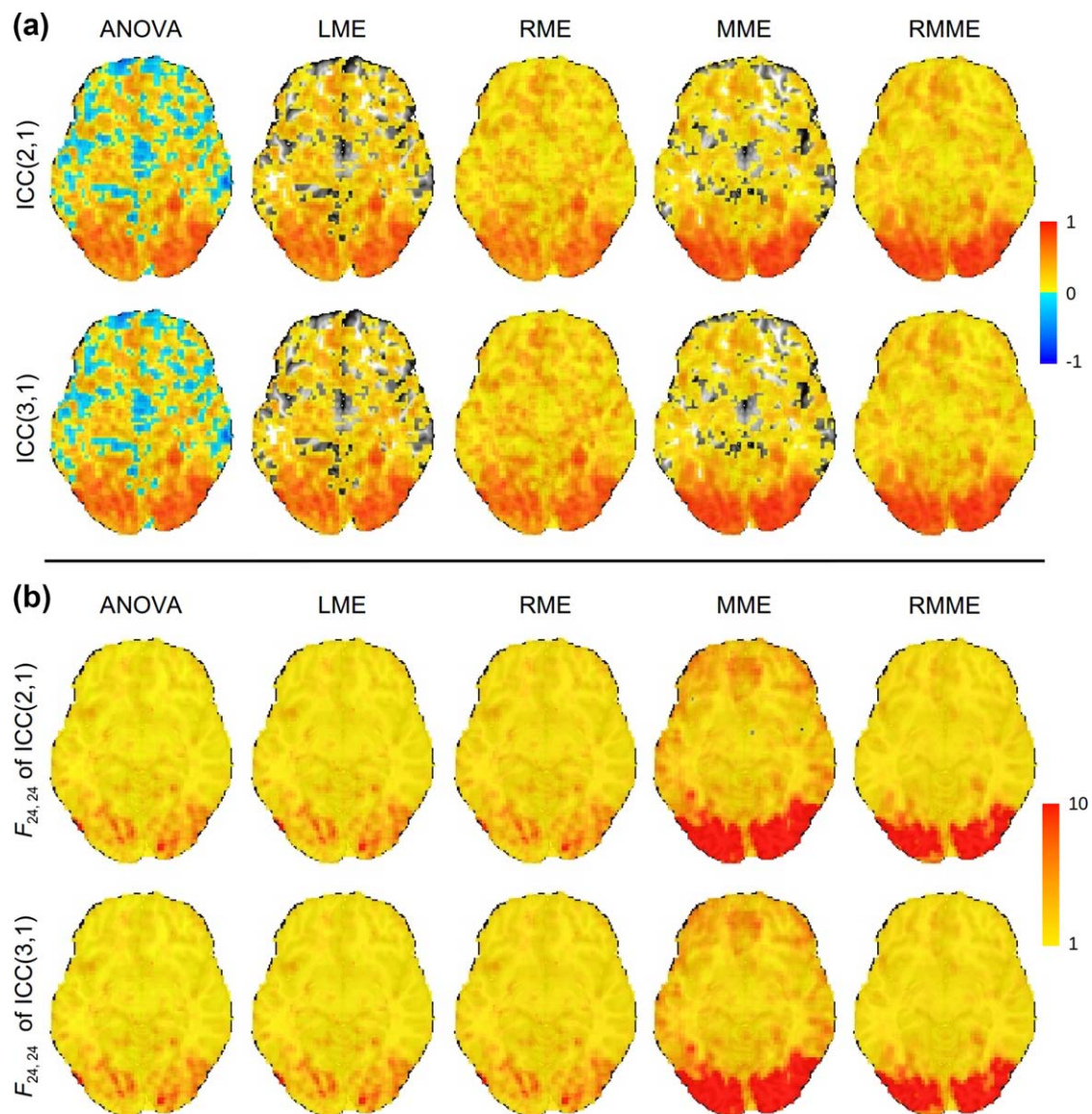


FIGURE 1 Panel A. ICC maps on an axial slice ($Z = -6$ mm, TT standard space; radiological convention: left is right) with a whole-brain dataset through four models. The conventional ANOVA (first column) and LME (second column) rendered substantial number of voxels with negative (blue) and zero (not colored) ICC values, respectively. Panel B. The F -statistic maps corresponding to each of the ICC maps in Panel A are shown with colors in the range of $[1, 10]$, with the lower bound of 1 defined by the formulation of F -statistic in (5). In each case, the degrees of freedom were $(24, 24)$ [Color figure can be viewed at wileyonlinelibrary.com]

0.5 (with η fixed at 2); their differences in ICC estimates across the four κ values are negligible. Furthermore, the computation cost is a decreasing function of κ (substantially highest at $\kappa = 0$). In light of these results, we set an empirical prior of gamma density (16) at $\eta=2, \kappa=0.5$ for neuroimaging data, and the choice is also consistent with the simulation results and recommendations in Chung et al. (2013).

For the voxels with negative ICC values from ANOVA or with zero ICC values from LME, RME offers positive but generally small ICC estimates; the nudging effect of RME is relatively small when the ICC value from ANOVA/LME is positive but small (<0.3), and it is negligible when the ICC value from ANOVA/LME is moderate or large (>0.3) (third column in Panel A, Figure 1; scatterplot cells (1, 3) and (3,1) in Figure 3). Some of the ICC estimates from MME are larger to varying

extent than ANOVA/LME/RME, while some are smaller (fourth column in Panel A in Figure 1; “fat blobs” showing wide variation in the scatterplot cells (1, 4), (2, 4), (3, 4), and their symmetric counterparts in Figure 3); there are a higher number of voxels with larger ICC from MME relate to ANOVA/LME/RME than those with smaller ICC values (slightly redder voxels in fourth column than the first three columns, Panel A in Figure 1; more dots with ICC greater than 0.75 on the MME side of the green diagonal line in the scatterplot cells (1, 4), (2, 4), (3, 4), and their symmetric counterparts in Figure 3). However, there are still a large fraction of voxels with zero ICC estimates from MME (uncolored voxels in the fourth column, Panel A, Figure 1), although less than from LME. Last, RMME shares similar properties with MME relative to ANOVA/LME/RME (fifth column, Panel A, Figure 1; scatterplot cells (1,

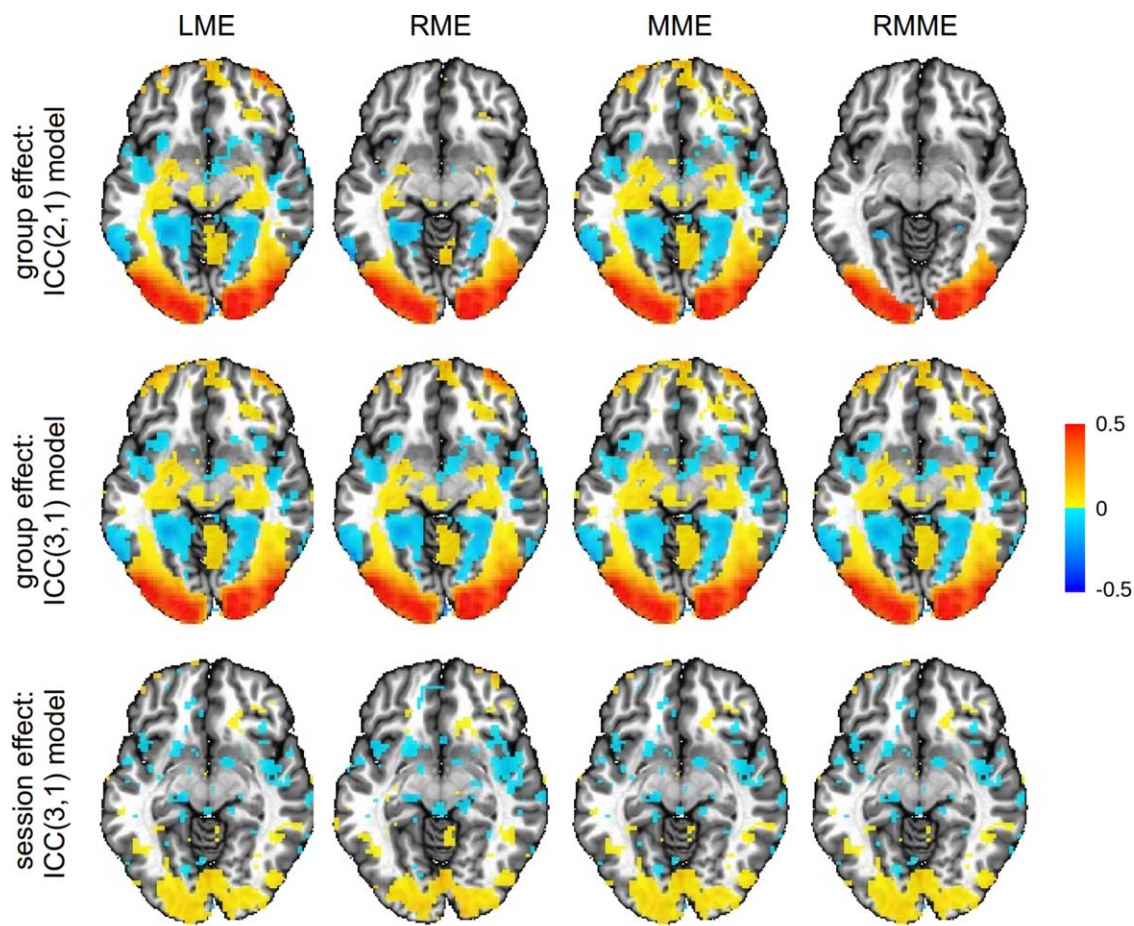


FIGURE 2 The first two rows show the group average effect on an axial slice ($Z = -6$ mm, TT standard space; radiological convention: left is right) with a whole-brain dataset estimated through the four mixed-effects models. The third row is the session effect in the ICC(3,1) model. While the maps display the session effect, they are thresholded at a liberal significance level of 0.1 with 24 degrees of freedom for the corresponding t -statistic. The estimates for the fixed effects are not directly available from ANOVA, and therefore are not displayed. As the two regularization approaches estimate slightly higher variances under the ICC(2,1) model, some regions in the frontal area fail to survive the liberal thresholding for RME and RMME as shown in the first row. In contrast, all the four models demonstrate similar results under the ICC(3,1) model as shown in the second row [Color figure can be viewed at wileyonlinelibrary.com]

5), (2, 5), (3, 5), and their symmetric counterparts Figure 3). The differences between RMME and MME parallel those between RME and LME (fifth column, Panel A, Figure 1; scatterplot cells (4, 5) and (5, 4), Figure 3); that is, RMME provides small but positive ICC estimates for those voxels with zero ICC under MME, and the nudging effect is negligible when the effect estimate is relatively reliable. The F -values across the five models follow roughly the same patterns as the ICC values (Panel B in Figure 1) because the F formulation is closely related to its ICC counterpart, and the high reliability measure in the brain shares the same regions with group average effect revealed from the four mixed-effects models (first two rows in Figure 2).

The differences between ICC(2,1) and ICC(3,1) (Panel A in Figure 1; diagonal cells in Figure 3) are mostly small for each of the five models except for those regions where session effect is substantial (third row, Figure 2). Because any systematic differences between the two sessions are accounted in the ICC(3,1), but not ICC(2,1), model, the former tends to render slightly higher ICC estimates. This is demonstrated in Figure 3, by the fact that most voxels are above the green diagonal line in the

diagonal cells of scatterplots. In addition, one noteworthy phenomenon is that RMME narrows the differences between the two ICC types, as represented in Figure 3 by the thinner band in scatterplot cell (5, 5) relative to the other diagonal cells. RMME tends to slightly overestimate ICC(3,1) relative to MME due to regularization as shown in the cell (5, 4); in contrast, such an upward pooling effect on ICC(2,1) relative to MME, as shown in the cell (4, 5), is slightly larger than ICC(3,1). The net impact of these two small upward pooling effects seems to bring the two ICC types close to each other, as shown in the cell (5, 5).

To gain better insights into the relative performances of the five models, we demonstrate three scenarios with three representative voxels⁹ with Table 4 illustrating the differences among the five models and with Figure 4 schematically showing the heterogeneity and heteroscedasticity at those three voxels (their effect estimates and their

⁹It might be neurologically more interesting to show the ICC in a few regions, but here we chose these three voxels instead of regions to demonstrate their subtle differences across ICC types as well as models.

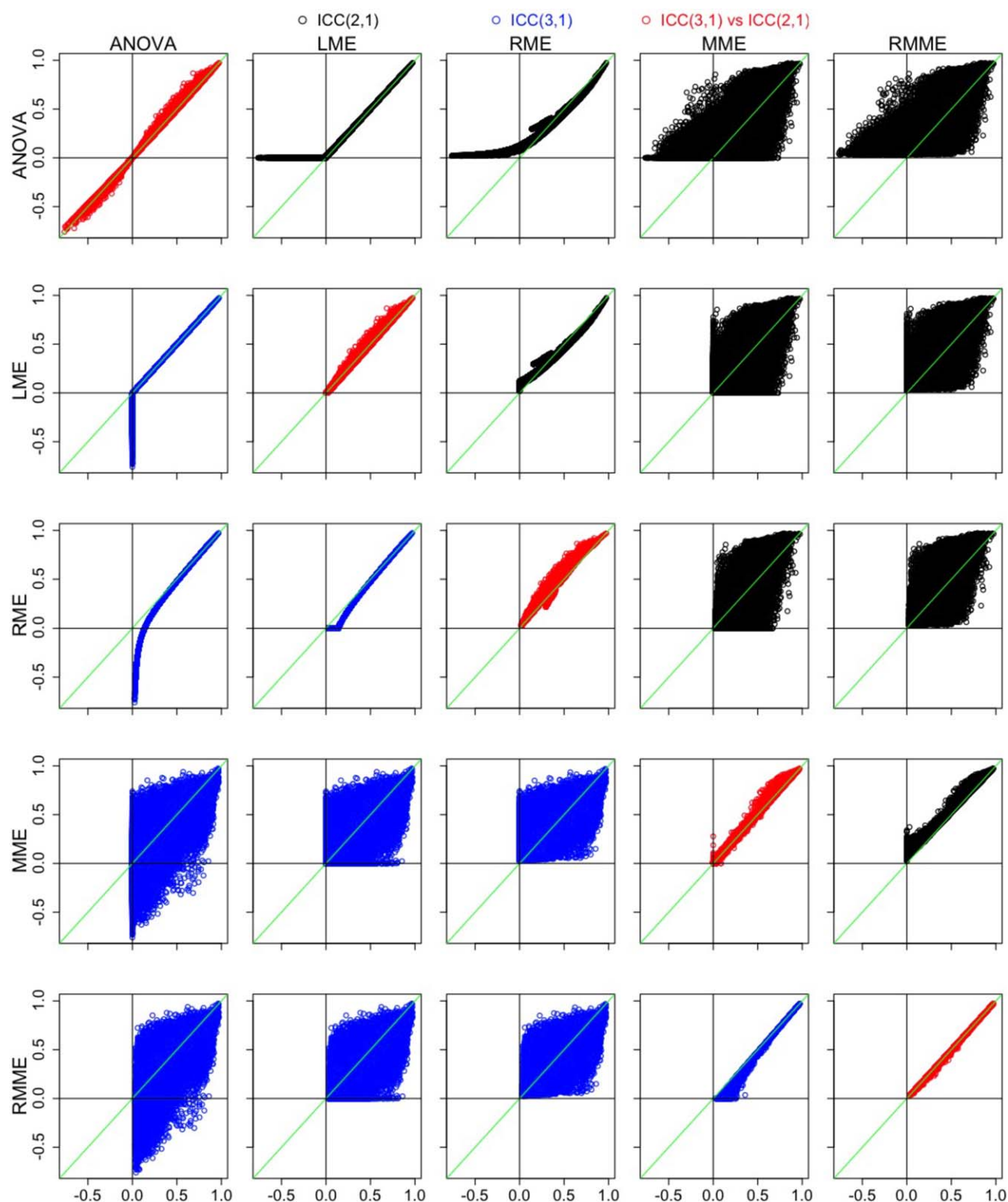


FIGURE 3 Comparisons across the five models and between ICC(2,1) and ICC(3,1) through scatterplots of ICC with the voxels in the brain. The x and y axes are a pair of ICC values between two of the five models with their labels shown at the left (x axis) and top (y axis), respectively. The 10 combinatorial comparisons for ICC(2,1) are illustrated in the upper triangular cells in black; the 10 combinatorial comparisons for ICC(3,1) are listed in the lower triangular cells in blue; and the diagonals compare the two ICC types in red for each of the five models with ICC(2,1) and ICC(3,1) as x and y axes, respectively. In each plot, the green diagonal line marks the situation when ICC(2,1) = ICC(3,1) [Color figure can be viewed at wileyonlinelibrary.com]

TABLE 4 Results at three example voxels^a shown in Figure 4

Voxel	Model	ICC(2,1)			ICC(3,1)			Session effect		
		Value	$F_{24,24}$	p value	Value	$F_{24,24}$	p value	Value	t_{24}	p value
V_1^*	ANOVA	0.530	3.300	0.0024	0.530	3.300	0.0024	-	-	-
	LME	0.531	3.292	0.0025	0.534	3.292	0.0025	0.012	1.144	0.26
	RME	0.500	3.578	0.0014	0.552	3.468	0.0017	0.012	1.159	0.26
	MME	0.504	3.033	0.0043	0.504	3.030	0.0043	0.008	0.786	0.44
	RMME	0.529	3.246	0.0027	0.527	3.231	0.0028	0.008	0.789	0.44
V_2	ANOVA	-0.270	0.560	0.920	-0.280	0.560	0.920	-	-	-
	LME	0	1	0.5	0	1	0.5	0.073	1.469	0.15
	RME	0.044	1.126	0.39	0.058	1.123	0.39	0.073	1.499	0.15
	MME	0.470	4.464	2.5e-4	0.631	4.422	2.7e-4	0.091	4.876	5.7e-5
	RMME	0.652	4.744	1.5e-4	0.649	4.693	1.7e-4	0.091	4.878	5.7e-5
V_3	ANOVA	0.570	4.800	1.4e-4	0.660	4.800	1.4e-4	-	-	-
	LME	0.574	4.819	1.4e-4	0.656	4.819	1.4e-4	0.080	3.372	0.0025
	RME	0.509	5.072	8.9e-5	0.665	4.970	1.1e-4	0.080	3.398	0.0023
	MME	0.743	13.851	5.8e-9	0.864	13.748	6.3e-9	0.076	5.52	1.1e-5
	RMME	0.870	14.435	3.8e-9	0.870	14.344	4.0e-9	0.076	5.517	1.1e-5

^aThe locations for the three voxels are the following (coordinates are in mm in TT space): V_1 : (39, 64, -6) in right middle occipital gyrus; V_2 : (11, 88, -11) in right lingual gyrus; V_3 : (11, 86, -6) in right lingual gyrus. Any fixed effects such as session difference here can be easily estimated through any of the four mixed-effects models for ICC(3,1), but they are usually not provided through ANOVA.

variances are listed in Appendix A). At voxel V_1 (from the right middle occipital gyrus), the ICC estimates from the ANOVA, LME and RME approaches are similar across all the five models. Both MME and RMME estimates are not much different because those effect estimates are roughly equally precise (first two columns in the right panel). The ICC estimate for the two types do not differ much either as the difference between the two sessions is small ($\sim 0.01\%$ signal change, Table 4).

At voxel V_2 (from the right lingual gyrus) ANOVA produces negative ICC values because $MS_\lambda < MS_\epsilon$ in the corresponding ICC formulation (Equation 2). LME avoids negativity by taking the lower boundary value of zero that is allowed for variance, while RME renders a positive but small ICC estimate. In contrast, both MME and RMME achieve more accurate variance estimates in the sense that the availability of variances for measurement errors provides a way to reallocate or redistribute the components in the total variance. In other words, the substantial amount of heteroscedasticity as shown in Figure 4 (third and fourth column in the right panel) allows differential weightings among the effect estimates with a higher emphasis on the more precise ones and downgrading for the less reliable ones. The dramatically different ICC values at voxel V_2 from MME and RMME, relative to the three other ICC models, can be understood by examining the wide range of effect estimates as well as their heterogeneous variances, as shown in the wide color spectrum of Figure 4. It is worth noting that the presence of substantial session effect (close to 0.1%, $p = 5.7 \times 10^{-5}$) leads to a moderate difference between ICC(2,1) and ICC(3,1) for MME, but the analogous difference for RMME is negligible; that is, RMME tends

to render similar ICC value between the two types regardless of the presence of session effect, just as shown in the whole-brain data (scatterplot cell (5, 5) in Figure 3).

Last, at voxel V_3 (from the right lingual gyrus), ANOVA, LME and RME reveal moderately reliable effect estimates with similar ICC estimates. However, both MME and RMME render higher ICC values, due to the presence of moderate amount of heteroscedasticity, similar to the situation with voxel V_2 even though less dramatic here (last two columns in both panels, Figure 4). Also similar to voxel V_2 , the session effect ($\sim 0.08\%$, $p = 1.1 \times 10^{-5}$) results in a lower estimate from MME for ICC(2,1) than ICC(3,1), but RMME estimates virtually the same reliability between the two ICC types. It is also interesting to note that, at both voxels V_2 and V_3 , the t -statistic for the fixed effect of session is much higher when the precision information is considered than in the other two models, LME and RME (Table 4). This phenomenon demonstrates the potential impact and importance of including modeling precision in neuroimaging group analysis (Chen et al., 2012; Worsley et al., 2002; Woolrich et al., 2004).

7 | DISCUSSION

Reliability is a crucial foundation for scientific investigation in general, and it has been a challenging issue and a hot topic for neuroimaging in particular over the years (Bennett & Miller, 2010). ICC offers a metric that can measure reliability under relatively strict circumstances. If the same scanning parameters are applied to the same cohort of subjects or families that undergo the same set of tasks, ICC shows the reliability

TABLE 5 A compact tcsh script that contains the succinct, selected *afni_proc.py* command used to generate the full processing pipeline (>500 lines) in AFNI for this study

```
#!/bin/tcsh
# Set top level directory structure
set subjID=$1
set currDir='pwd'
set anatDir=./freesurfer.anat/${subjID}/mri
set epiDir=afni
set stimDir=stim.files
# miscellaneous parameters or options
set motion_max=1.0; set delete_nfirst=4; set costfunc=lpc+zz
# run afni_proc.py to create a single subject processing script
afni_proc.py -subj_id ${subjID} \
-script proc.script.FINAL.1.20.17.${subjID} -scr_overwrite\
-blocks despike tshift align tlrc volreg blur mask scale regress\
-copy_anat $anatDir/brainmask.nii\
-tcat_remove_first_trs $delete_nfirst\
-dsets $epiDir/r01+orig $epiDir/r02+orig $epiDir/r03+orig\
-blur_size 5 -out_dir NL.results\
-anat_unif_GM no -anat_has_skull no -tlrc_NL_warp\
-volreg_align_e2a -volreg_align_to MIN_OUTLIER -volreg_tlrc_warp\
-align_opts_aea -cost $costfunc -giant_move -AddEdge\
-regress_stim_times\
$stimDir/${subjID} n_A_timing.1D\
$stimDir/${subjID} a50_timing.1D\
$stimDir/${subjID} a100_timing.1D\
$stimDir/${subjID} a150_timing.1D\
$stimDir/${subjID} n_F_timing.1D\
$stimDir/${subjID} f50_timing.1D\
$stimDir/${subjID} f100_timing.1D\
$stimDir/${subjID} f150_timing.1D\
$stimDir/${subjID} n_H_timing.1D\
$stimDir/${subjID} h50_timing.1D\
$stimDir/${subjID} h100_timing.1D\
$stimDir/${subjID} h150_timing.1D\
$stimDir/${subjID} w_timing.1D\
-regress_stim_labels\
n_A a50 a100 a150 n_F f50 f100 f150 n_H h50 h100 h150 w\
-regress_local_times -regress_censor_outliers 0.1\
-regress_basis 'BLOCK(2,1)' -regress_censor_motion $motion_max\
-regress_est_blur_epits -regress_est_blur_errts\
-regress_reml_exec -regress_compute_fitts -regress_opts_3d\
-allzero_OK -regress_opts_reml -GOFORIT\
-regress_make_ideal_sum sum_ideal.1D\
-gltsym 'SYM: +n_A +a50 +a100 +a150 +n_F +f50 +f100 +f150 +n_H +h50 +h100 +h150'\
-glt_label 1 Positive_Control
tcsh -xef proc.script.${subjID} |& tee proc.script.${subjID}.output
```

Note. To implement across the group, one simply loops through a list of subjects, entering the given file name as the sole command line argument, which is passed to the variable \$subjID. Here, stimulus variables are encoded as: h = happy, n = neutral, f = fearful, w = wrong; and each is followed by the duration (50, 100, 150 s).

across the levels of a categorical variable such as twins, sessions, or scanners. Nevertheless, in practice it is difficult to keep scanning circumstances perfectly constant across sessions in neuroimaging, thus systematic differences may slip in one way or another. Additionally, ICC can be applied to neuroimaging in the classic sense (i.e., inter-rater reliability or concordance), assessing the reliability between two diagnostic approaches (e.g., human versus automatic method) on a disease (e.g., schizophrenia, autism, or depression) or two different analytical methods applied to the same collection of data.

As shown here, there are three types of ICC, each of which can be estimated through various models such as ANOVA, LME, RME, MME, and RMME. On one hand, the ICC metric offers a unique approach to measuring the reliability of neuroimaging data under some well-controlled conditions; on the other hand, the investigator may still face a daunting job in deciding which ICC type, and which model, is most

appropriate to apply in a given experiment or setup. We have presented the different formulations of each here, and demonstrated differences in data outcomes. We further discuss and summarize recommendations for model selections below.

7.1 | Considerations for effect estimates as inputs for ICC analysis

In practice, several factors can contribute to having poor data quality and accuracy in neuroimaging. Having relatively low signal-to-noise ratio is a major issue, and suboptimal modeling due to poor understanding of the major components in the signal is another. By some estimates, less than half (and in some cases down to 20%–30%) of data variability can be accounted for in a typical fMRI data analysis at the individual level (Gonzalez-Castillo et al., 2017). Some rigor and standardization steps are

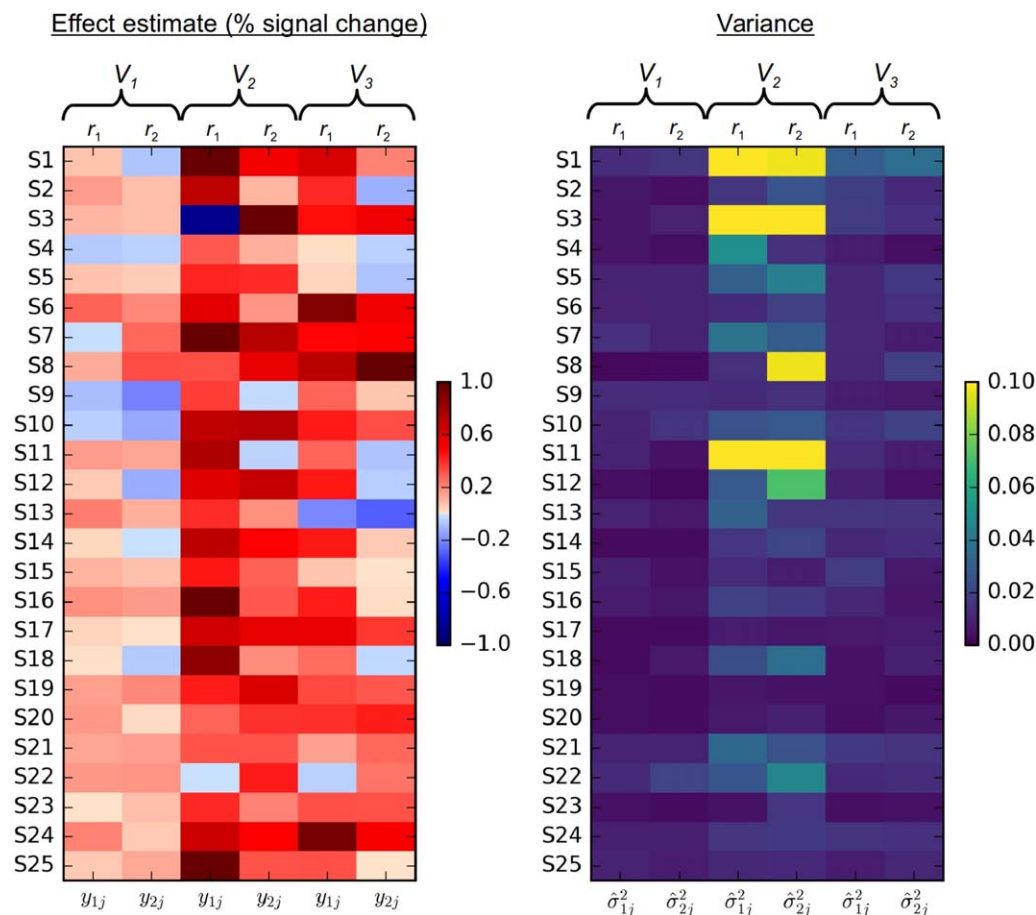


FIGURE 4 Example illustrations of heterogeneity and heteroscedasticity at three voxels with results shown in Table 4. Effect estimates (left) during the two sessions (r_1 and r_2) and the corresponding variances (right) for the measurement errors for the 25 subjects at three gray matter voxels (V_1 , V_2 , and V_3) from the axial slice shown in Figure 1. The real data are listed in Appendix A. Substantial session effect can be seen in voxel V_2 (left matrix, middle two columns) while the session effect is negligible at voxels V_1 and V_3 . A large amount of variability exists for the effect estimates at voxel V_1 (left matrix, first two columns), leading to negative ICC estimate by ANOVA and zero by LME, while RME manages to deal with the degenerative situation with a small, but positive, ICC estimate. MME provides a positive and relatively large ICC estimate through weighting based on the precision information (right matrix, first two columns). The variability for the precision of the effect estimate is moderate at voxel V_2 (right matrix, middle two columns), but minimal at voxel V_3 (right matrix, right two columns) [Color figure can be viewed at wileyonlinelibrary.com]

required to achieve more accurate reproducibility and reliability. At present, the cumulative impact of altering preprocessing and modeling strategy in an analysis pipeline is largely unknown.

For example, it is well known that the absolute values of the fMRI-BOLD signal from the scanner have arbitrary units with scaling fluctuations among subjects, and therefore some kind of calibration is needed during preprocessing if the effect estimates are to be used in further analyses at the group level, for both typical group analysis and ICC estimation. Such a calibration should take into consideration the fact that the baseline varies across brain regions, subjects, groups, scanners, and sites. However, radically different scaling strategies are in use, and their effects at the group level remain unclear. For instance, the global or grand mean scaling, as currently typically practiced and adopted by some software packages, can only account for part of the cross-subject variability; however, such approaches do not address the cross-region variability within an individual brain as a practical reality in neuroimaging scanning, and therefore, either may lead to difficulty in interpreting and comparing the effect estimates. Analyses on networks including

ICC and causal modeling would be affected as well. In contrast, voxel-wise scaling as a calibrator, even though imperfect due to the ambiguity in the baseline definition (Chen, Taylor, & Cox, 2017b), provides a more accurate characterization than the alternatives. Because of these considerations, in general, we recommend voxel-wise mean scaling during preprocessing, so that a directly interpretable and compatible measure in percent signal change can be adopted for the effect estimates that are to be taken for group analysis including ICC.

Last, each effect estimate in task-related experiment is typically derived from a model coefficient, and is intrinsically associated with some extent of uncertainty that may vary across subjects. As the precision information of the single subject effect estimate (embedded in the denominator of t-statistic) is required for the preferred MME and RMME, it is important to model the temporal correlation in the residuals at the individual level to avoid inflating the precision (or underestimating the uncertainty) for the effect estimate to be used in the ICC model. As for resting-state data, the precision information for the correlation coefficient is not directly available. However, it is known that

the Fisher-transformed Z-value for the correlation coefficient between any two time series that are white noise approximately follows a Gaussian distribution $N(\frac{1}{2} \ln \frac{1+r}{1-r}, \frac{1}{T-3})$ (Sheskin, 2004), where r is the Pearson correlation of two time series each having T time points. Therefore, $\frac{1}{T-3}$ is the lower bound for the variance of correlation coefficient in the brain. Nevertheless, the variance is usually not expected to have a substantial variability across the brain and across subjects, and thus the lack of precision information is not considered an issue for ICC computation with resting-state data analysis.

7.2 | Which ICC model to use?

Among the thousands of voxels in a typical whole-brain neuroimaging dataset, negative ICC values unavoidably, and even frequently, show up, though they are usually not reported in the literature. The standard ICC values reported in most literature to date contain several aspects of ambiguity, greatly hindering meaningful interpretation. In addition, even when the ICC type is reported, its reason for selection, for example, between ICC(2,1) and ICC(3,1), is usually not clearly explained. In some cases, the chosen method is ill-suited to the given analysis scenario (e.g., if ICC(1,1) was used between two sessions).

When precision information for the measurement errors of the effect estimate is available, we recommend using RMME for the following three reasons: (a) the precision information offers a more robust estimate for ICC, and for the fixed effects in the model; (b) the regularization aspect of the approach leads to the avoidance of the uninterpretable situation of a negative ICC from ANOVA or an unrealistic zero ICC estimate from plain LME; (c) as demonstrated here (scatterplot cell (5, 5) in Figure 3 and Table 4), RMME tends to be less sensitive to ICC type selection, rendering roughly the same ICC estimate regardless of the type the investigator adopts. When the precision information is unavailable (e.g., when one has the correlation value as the effect of interest from seed-based analysis or psycho-physiological interaction analysis), we recommend using RME because of its capability to provide a realistic ICC estimate when ANOVA (or LME) renders negative (or zero) ICC.

A general linear model (GLM), as extension to ANOVA, can accommodate between-subjects variables and assist the investigator in ICC type selection. Therefore, the GLM is an intermediate approach between ANOVA and LME. However, it still shares the other limitations with ANOVA in the following aspects: rendering negative or zero ICC, and being unable to handle missing data or to incorporate sampling errors.

The use of a regularization approach may raise questions regarding its insertion of arbitrariness into the estimation process with a prior for the variance components. Each variance component in the ICC formulation is estimated as a point value; however, it is worth noting that the value of a variance component usually does not exactly fall at its numerical estimate, but varies within some range (e.g., 95% central or uncertainty interval). The reason for a negative or zero estimate for a variance component lies in the methodology of replacing the variance by a point estimate; in other words, the substitution with a point estimate ignores the fact that there is uncertainty associated with the estimate. As the

point estimate usually tends to be imprecise and underestimated, the negative or zero variance estimate or ICC should not be taken at the face value (Chung et al., 2013). A zero estimate for variance and ICC can also cause inflated inferences on the fixed effects such as group average as well as systematic differences across the factor A levels in the ICC (3,1) model. More fundamentally, simply forcing a negative ICC or the variance estimate for σ_λ^2 in Equation 2 to zero leads to an unjustifiable claim that the effects from all the subjects are absolutely the same. On the other hand, the regularization approach can be conceptualized as a tug of war between the prior and data. As the results from our experimental dataset demonstrate here (e.g., Figure 1), a weakly informative prior for ICC estimation can pull a degenerate situation (e.g., zero variance estimate) out of the boundary and render a more reasonable estimate, while little impact will incur when the data contain strong information, which would overrule the prior.

Four aspects of MME and RMME are worth highlighting here. First, these two multilevel approaches estimate ICC differently from the other three methods to some extent, as indicated in those “fat blobs” of Figure 3 among cells (1,4), (2,4), (3,4), (1,5), (2,5), and (3,5) for ICC(2,1), and (4,1), (4,2), (4,3), (5,1), (5,2), and (5,3) for ICC(3,1). It should be stressed that the focus here is not on which method leads to a higher or lower estimate (nor should it be), but on which model provides a more accurate characterization about the reliability measure. Second, just as the levels of a within-subject factor are treated as simultaneous variables of a multivariate model (as in the AFNI group analysis program 3dMVM) for repeated-measures ANOVA in neuroimaging group analysis¹⁰ (Chen et al., 2014), so are the factor A levels in MME and RMME for ICC(3,1) modeled here in a multivariate fashion with the flexibility in variance decomposition (Viechtbauer, 2010) and the capability to incorporate quantitative covariates in the presence of a repeated-measures factor. Third, the measurement errors associated with the factor A levels in a generic model of the form in Equation 1 can be correlated when different tasks are intertwined in the experiment, thus the variance-covariance matrix \mathbf{R} should be semi-definite in general. It is because of the presence of covariances in \mathbf{R} that the modeling approach adopted in 3dMEMA of AFNI and FLAME of FSL cannot instead be utilized to perform a paired test in cases where the two effect estimates are entered separately as input (even if the program permits such an option), as the measurement errors corresponding to the two effect estimates are usually correlated.¹¹ However, \mathbf{R} is actually a diagonal matrix in the neuroimaging ICC context, because the measurement errors can be reasonably assumed to be well-approximated as independent with each other (e.g., across runs or sessions). Finally, nonparametric methods (e.g., bootstrapping, permutations) may offer a robust venue for controlling family-wise error for relatively simple models at the group level; however, under some scenarios, parametric approaches provide more specific and accurate characterization about the effect of interest

¹⁰Due to the complexity of handling within-subject factors, some group analyses in published Neuroimaging literature are still not performed correctly, leading to substantial number of publications with inflated statistical inferences (McLaren, 2010; Chen et al., 2014).

¹¹Instead, such a paired test can be validly performed by taking the contrast and its standard error as input.

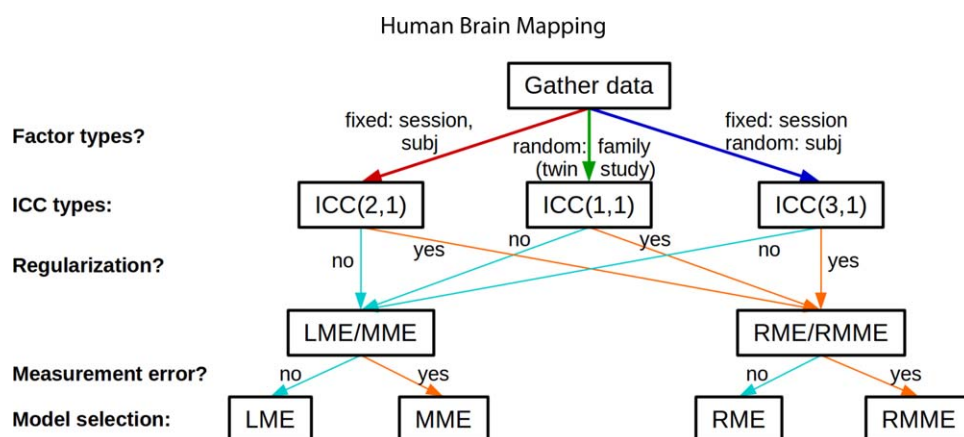


FIGURE 5 Flowchart of ICC modeling options [Color figure can be viewed at wileyonlinelibrary.com]

through some quantifying parameters (e.g., variance components in the ICC context) in the model, which are currently both valid and irreplaceable, as shown here with ICC computations through ANOVA, LME, RME, MME, and RMME, and the group analysis approach through incorporating the precision information (Chen et al., 2012; Woolrich et al., 2004; Worsley et al., 2002).

7.3 | Which ICC type to adopt?

As there is usually one single effect estimate for each subject per scanning situation, our discussion here focuses on single-measurement ICC. Among the three ICC types, ICC(1,1) is likely the easiest to list the scenarios in which it can be applied. Primarily, it can be used when there is no apparent distinction for a sequence of the levels for the factor that is associated with the multiple measurements in the ICC model. It typically applies to the situation of studying twins, for example.

In contrast, the other two types are utilized for scenarios such as having two or more runs, sessions, scanners, sites, or between parent and child. The reliability from ICC(2,1) represents “absolute agreement,” to the extent that the values exactly match between any two levels of the factor A, while ICC(3,1) shows the consistency or the extent that the effect estimates match across the factor A levels *after* accounting for potential systematic differences or other fixed effects. In other words, if the systematic differences across the levels of factor A are negligible, then the two ICC estimates would be similar. On the other hand, if the systematic differences or confounding effects are substantial, then the ICC values tend to diverge to some extent, and they lead to different interpretations. However, the existence of systematic differences itself warrants further exploration about the source or nature of those fixed effects (e.g., habituation or attenuation). For example, what is the association between the ICC map and the activation map (i.e., intercept in the model (Equation 14))?

Owing to the simultaneity of analyzing all the voxels, it is unrealistic to choose one ICC type for some regions while selecting another for the rest. Per the discussion here between ICC(2,1) and ICC(3,1), we generally recommend the latter for whole-brain analysis. In doing so, potential fixed effects are properly accounted for. More importantly, it is not the ICC interpretation in the sense of absolute agreement that is

generally of primary importance, but the extent of agreement after potential fixed effects are all explained. Furthermore, with ICC(3,1), the investigator can directly address the following questions: (a) Which brain regions show systematic effects across the levels of factor A? (b) How do those systematic effects correspond to the ICC maps in these regions? (c) Are the systematic effects related to some confounding effects such as habituation or attenuation?

The overall decision tree for ICC computation is summarized as a flowchart in Figure 5.

7.4 | Result reporting

Clear and accurate scientific reporting is important for result reproducibility, reliability and validity as well as for meta-analysis. The present reporting conventions in neuroimaging are especially discouraging, with a lopsided focus on statistics alone, for example, due to oversight or limitations in software implementations (Chen et al., 2017b), leading to incomplete reporting through the literature. It cannot be emphasized enough that the effect estimates involved in a study should be reported in addition to the statistical significance, and the same emphasis should be applied to ICC analysis. Specifically, the investigator should explicitly state the ICC type and the model adopted, and the justification for such choices. One may notice that the ICC formulation for ICC(1,1) in Equations 11 and 12 is exactly the same as ICC(3,1) in Equations 7 and 8, which means that reporting the ICC formula would not be enough to reveal the whole picture because the two underlying models (Equations 10 and 6) are dramatically different (and so are the two resulting ICC estimates).

With regards to the criteria for reliability, a loose rule of thumb has been suggested for ICC values as the following (Cicchetti, 1994): [0, 0.4), poor; [0.4, 0.6), fair; [0.6, 0.75), good; and [0.75, 1], strong. One cautionary note is that a low ICC does not always mean poor reliability: it is possible that some confounding effects are not accounted for in the model. For the statistical significance of ICC, one may use the Fisher transformation (Equation 4), or preferably, the *F*-statistic (Equations 5 and 9). Nevertheless, the *F*-statistic is not necessarily a basis for clusterization, but together with the ICC value, it serves as some auxiliary information to gauge the reliability about the effect of interest from the conventional

analytical pipeline. Finally, all the fixed effects including the intercept (group average) are crucial part of the model and should be discussed and explained in the article as well, as exemplified here in Figures 1 and 2.

8 | CONCLUSION

One potential problem with the classic definition of ICC is that negative ICC values may appear a scenario which is almost certain to show up in a whole-brain neuroimaging analysis. Here we extend the conventional ICC definition and its computations to the frameworks of LME, RME, MME, and RMME modeling to address this issue and other difficulties. Such an extension not only offers wider modeling flexibility such as the inclusion of various fixed effects but also avoids the interpretability problem of negative ICC values under ANOVA. We offer our recommendations in model adoption and ICC type selection and also in result reporting. All modeling strategies and ICC types and the estimation and statistic testing for the fixed effects are currently available in the AFNI program 3dICC.

ACKNOWLEDGMENTS

The research and writing of the article were supported by the NIMH and NINDS Intramural Research Programs (ZICMH002888) of the NIH/HHS, USA. The authors are indebted to Wolfgang Viechtbauer and Vincent Dorie for their precious help and for the R packages *metafor* and *brms*, which they build and maintain, respectively.

ORCID

Gang Chen  <http://orcid.org/0000-0002-2960-089X>

REFERENCES

- Baayen, R., Davidson, D., & Bates, D. (2008). Mixed-effects modeling with crossed random effects for subjects and items. *Journal of Memory and Language*, 59, 390–412.
- Bates, D. M., Maechler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1), 1–48.
- Bennett, C. M., & Miller, M. B. (2010). How reliable are the results from functional magnetic resonance imaging? *Annals of the New York Academy of Sciences*, 1191, 133–155.
- Bhambhani, Y., Maikala, R., Farag, M., & Rowland, G. (2006). Reliability of near-infrared spectroscopy measures of cerebral oxygenation and blood volume during handgrip exercise in nondisabled and traumatic brain-injured subjects. *The Journal of Rehabilitation Research and Development*, 43(7), 845–856.
- Braun, U., Plichta, M. M., Esslinger, C., Sauer, C., Haddad, L., Grimm, O., ... Meyer-Lindenberg, A. (2012). Test-retest reliability of resting-state connectivity network characteristics using fMRI and graph theoretical measures. *NeuroImage*, 59(2), 1404–1412.
- Brandt, D. J., Sommer, J., Krach, S., Bedenbender, J., Kircher, T., Paulus, F. M., & Jansen, A. (2013). Test-retest reliability of fMRI brain activity during memory encoding. *Frontiers in Psychiatry*, 4, 163.
- Cao, H., Plichta, M. M., Schfer, A., Haddad, L., Grimm, O., Schneider, M., ... Tost, H. (2014). Test-retest reliability of fMRI-based graph theoretical properties during working memory, emotion processing, and resting state. *NeuroImage*, 84, 888–900.
- Cáceres, A., Hall, D. L., Zelaya, F. O., Williams, S. C. R., & Mehta, M. A. (2009). Measuring fMRI reliability with the intra-class correlation coefficient. *NeuroImage*, 45(3), 758–768.
- Chen, G., Adelman, N. E., Saad, Z. S., Leibenluft, E., & Cox, R. W. (2014). Applications of Multivariate Modeling to Neuroimaging Group Analysis: A Comprehensive Alternative to Univariate General Linear Model. *NeuroImage*, 99, 571–588.
- Chen, B., Xu, T., Zhou, C., Wang, L., Yang, N., Wang, Z., ... Weng, X.-C. (2015). Individual variability and test-retest reliability revealed by ten repeated resting state brain scans over one month. *PLoS ONE*, 10(12), e0144963.
- Chen, G., Saad, Z. S., Nath, A. R., Beauchamp, M. S., & Cox, R. W. (2012). FMRI group analysis combining effect estimates and their variances. *NeuroImage*, 60, 747–765.
- Chen, G., Saad, Z. S., Britton, J. C., Pine, D. S., & Cox, R. W. (2013). Linear mixed-effects modeling approach to FMRI group analysis. *NeuroImage*, 73, 176–190.
- Chen, G., Taylor, P. A., Shin, Y. W., Reynolds, R. C., & Cox, R. W. (2017a). Untangling the relatedness among correlations, part II: Inter-subject correlation group analysis through linear mixed-effects modeling. *NeuroImage*, 147, 825–840.
- Chen, G., Taylor, P. A., & Cox, R. W. (2017b). Is the statistic value all we should care about in neuroimaging? *NeuroImage*, 147, 952–959.
- Chung, Y., Rabe-Hesketh, S., Dorie, V., Gelman, A., & Liu, J. (2013). A nondegenerate penalized likelihood estimator for variance parameters in multilevel models. *Psychometrika*, 78(4), 685–709.
- Cicchetti, D. V. (1994). Guidelines, criteria, and rules of thumb for evaluating normed and standardized assessment instruments in psychology. *Psychological Assessment*, 6(4), 284–290.
- Cox, R. W. (1996). AFNI: Software for analysis and visualization of functional magnetic resonance neuroimages. *Computers and Biomedical Research*, 29, 162–173. <http://afni.nimh.nih.gov>.
- Fiecas, M., Ombao, H., van Lunen, D., Baumgartner, R., Coimbra, & Feng, D. (2013). Quantifying temporal correlations: A test-retest evaluation of functional connectivity in resting-state fMRI. *NeuroImage*, 65, 231–241.
- Fournier, J. C., Chase, H. W., Almeida, J., & Phillips, M. L. (2014). Model specification and the reliability of fMRI results: Implications for longitudinal neuroimaging studies in psychiatry. *PLoS ONE*, 9(8), e105169.
- Gonzalez-Castillo, J., Chen, G., Nichols, T., & Bandettini, P. A. (2017). Variance decomposition for single-subject task-based fMRI activity estimates across many sessions. *NeuroImage*, 154, 206–218.
- Griffanti, L., Rolinski, M., Szewczyk-Krolikowski, K., Menke, R. A., Filippini, N., Zamboni, G., Jenkinson, M., Hu, M. T. M., Mackay, C. E. (2016). Challenges in the reproducibility of clinical studies with resting state fMRI: An example in early Parkinson's disease. *NeuroImage*, 124, 704–713.
- Guo, C. C., Kurth, F., Zhou, J., Mayer, E. A., Eickhoff, S. B., Kramer, J. H., & Seeley, W. W. (2012). One-year test-retest reliability of intrinsic connectivity network fMRI in older adults. *NeuroImage*, 61(4), 1471–1483.
- Haller, S. P., Kircanski, K., Stoddard, J., White, L., Chen, G., Sharif-Askary, B., ... Brotman, M. A. (2017). Reliability of neural activation and connectivity during implicit face emotion processing in youth. In preparation.
- Higgins, J. P. T., & Thompson, S. G. (2002). Quantifying heterogeneity in a meta-analysis. *Statistics in Medicine*, 21(11), 1539–1558.
- Jaeger, L., Marchal-Crespo, L., Wolf, P., Riener, R., Kollias, S., & Michels, L. (2015). Test-retest reliability of fMRI experiments during robot-assisted active and passive stepping. *Journal of NeuroEngineering and Rehabilitation*, 12, 102.
- Kristo, G., Rutten, G.-J., Raemaekers, M., de Gelder, B., Rombouts, S. A. R. B., & Ramsey, N. F. (2014). Task and task-free FMRI

- reproducibility comparison for motor network identification. *Human Brain Mapping*, 35, 340–352.
- Lin, Q., Dai, Z., Xia, M., Han, Z., Huang, R., Gong, G., ... He, Y. (2015). A connectivity-based test-retest dataset of multi-modal magnetic resonance imaging in young healthy adults. *Scientific Data*, 2.
- McGraw, K. O., & Wong, S. P. (1996). Forming inferences about some intraclass correlation coefficients. *Psychological Methods*, 1, 30–46.
- Molloy, E. K., & Birn, R. M. (2014). Tools for test-retest fMRI studies. *ZENODO*, <https://doi.org/10.5281/zenodo.49326>.
- Noble, S., Scheinost, D., Finn, E. S., Shen, X., Papademetris, X., Mc Ewen, S. C., ... Constable, R. T. (2017). Multisite reliability of MR-based functional connectivity. *NeuroImage*, 146, 959–970.
- Pinheiro, J. C., & Bates, D. M. (2004). *Mixed-Effects Models in S and S-PLUS*. Springer.
- Plichta, M. M., Herrmann, M. J., Baehne, C. G., Ehlis, A. C., Richter, M. M., Pauli, P., & Fallgatter, A. J. (2006). Event-related functional near-infrared spectroscopy (fNIRS): Are the measurements reliable? *NeuroImage*, 31(1), 116–124.
- Plichta, M. M., Herrmann, M. J., Baehne, C. G., Ehlis, A. C., Richter, M. M., Pauli, P., & Fallgatter, A. J. (2007). Event-related functional near-infrared spectroscopy (fNIRS) based on craniocerebral correlations: Reproducibility of activation? *Human Brain Mapping*, 28(8), 733–741.
- Quiton, R. L., Keaser, M. L., Zhuo, J., Gullapalli, R. P., & Greenspan, J. D. (2014). Intersession reliability of fMRI activation for heat pain and motor tasks. *NeuroImage. Clinical*, 5, 309–321.
- Recasens, M., & Uhlhaas, P. J. (2017). Test-retest reliability of the magnetic mismatch negativity response to sound duration and omission deviants. *NeuroImage*, 157, 184–195.
- Revelle, W. (2016). *PSYCH: Procedures for Personality and Psychological Research*, Northwestern University, Evanston, Illinois, USA. <https://CRAN.R-project.org/package=psych> Version = 1.6.9.
- Shah, L. M., Cramer, J. A., Ferguson, M. A., Birn, R. M., & Anderson, J. S. (2016). Reliability and reproducibility of individual differences in functional connectivity acquired during task and resting state. *Brain and Behavior*, 6(5), e00456.
- Sheskin, D. J. (2004). *Parametric and nonparametric statistical procedures* (3rd ed.). Chapman & Hall/CRC.
- Shou, H., Eloyan, A., Lee, S., Zipunnikov, V., Crainiceanu, A. N., Nebel, M. B., ... Crainiceanu, C. M. (2013). Quantifying the reliability of image replication studies: The image intraclass correlation coefficient (I2C2). *Cognitive, Affective, & Behavioral Neuroscience*, 13(4), 714–724.
- Shrout, P., & Fleiss, J. (1979). Intraclass correlations: Uses in assessing rater reliability. *Psychological Bulletin*, 86(2), 420–428.
- Tian, F., Kozel, F. A., Yennu, A., Croarkin, P. E., McClintock, S. M., Mapes, K. S., ... Liu, H. (2012). Test-retest assessment of cortical activation induced by repetitive transcranial magnetic stimulation with brain atlas-guided optical topography. *Journal of Biomedical Optics*, 17(11), 116020.
- Tierney, L., Rossini, A. J., Li, N., & Sevcikova, H. (2016). SNOW: Simple Network of Workstations. R package version 0.4-2. <https://CRAN.R-project.org/package=snow>
- Töger, J., Sorensen, T., Somandepalli, K., Toutios, A., Lingala, S. G., Narayanan, S., & Nayak, K. (2017). Test-retest repeatability of human speech biomarkers from static and real-time dynamic magnetic resonance imaging. *Journal of the Acoustical Society of America*, 141, 3323.
- Viechtbauer, W. (2010). Conducting meta-analyses in R with the metafor package. *Journal of Statistical Software*, 36(3), 1–48.
- Wang, J. H., Zuo, X. N., Gohel, S., Milham, M. P., Biswal, B. B., & He, Y. (2011). Graph theoretical analysis of functional brain networks: Test-retest evaluation on short- and long-term resting-state functional MRI data. *PLoS ONE*, 6(7), e21976.
- White, L. K., Britton, J. C., Sequeira, S., Ronkin, E. G., Chen, G., Bar-Haim, Y., ... Pine, D. S. (2016). Behavioral and neural stability of attention bias to threat in healthy adolescents. *NeuroImage*, 136, 84–93.
- Woolrich, M. W., Behrens, T. E. J., Beckmann, C. F., Jenkinson, M., & Smith, S. M. (2004). Multilevel linear modeling for fMRI group analysis using Bayesian inference. *NeuroImage*, 21(4), 1732–1747.
- Worsley, K. J., Liao, C., Aston, J., Petre, V., Duncan, G. H., Morales, F., & Evans, A. C. (2002). A general statistical analysis for fMRI data. *NeuroImage*, 15, 1–15.
- Yan, C. G., Wang, X. D., Zuo, X. N., & Zang, Y. F. (2016). DPABI: Data processing and analysis for (resting-state) brain imaging. *Neuroinformatics*, 14, 339–351.
- Zanto, T. P., Pa, J., & Gazzaley, A. (2014). Reliability measures of functional magnetic resonance imaging in a longitudinal evaluation of mild cognitive impairment. *NeuroImage*, 84, 443–452.
- Zhang, H., Duan, L., Zhang, Y. J., Lu, C. M., Liu, H., & Zhu, C. Z. (2011). Test-retest assessment of independent component analysis-derived resting-state functional connectivity based on functional near-infrared spectroscopy. *NeuroImage*, 55(2), 607–615.
- Zuo, X. N., Martino, A. D., Kelly, C., Shehzad, Z. E., Gee, D. G., Klein, D. F., ... Milham, M. P. (2010a). The oscillating brain: Complex and reliable. *NeuroImage*, 49(2), 1432–1445.
- Zuo, X. N., Kelly, C., Adelstein, J. S., Klein, D. F., Castellanos, F. X., & Milham, M. P. (2010). Reliable intrinsic connectivity networks: Test-retest evaluation using ICA and dual regression approach. *NeuroImage*, 49(3), 2163–2177.
- Zuo, X. N., Xu, T., Jiang, L., Yang, Z., Cao, X. Y., He, Y., ... Milham, M. P. (2013). Toward reliable characterization of functional homogeneity in the human brain: Preprocessing, scan duration, imaging resolution and computational space. *NeuroImage*, 65, 374–386.
- Zuo, X. N., & Xing, X. X. (2014). Test-retest reliabilities of resting-state FMRI measurements in human brain functional connectomics: A systems neuroscience perspective. *Neuroscience & Biobehavioral Reviews*, 45, 100–118.
- Zuo, X. N., Anderson, J. S., Bellec, P., Birn, R. M., Biswal, B. B., Blautzik, R. L., ... Milham, M. P. (2014). An open science resource for establishing reliability and reproducibility in functional connectomics. *Scientific Data*, 1, 140049.

How to cite this article: Chen G, Taylor PA, Haller SP, et al. Intraclass correlation: Improved modeling approaches and applications for neuroimaging. *Hum Brain Mapp*. 2018;39:1187–1206. <https://doi.org/10.1002/hbm.23909>

APPENDIX A

FMRI PROCESSING

The general sequence of fMRI data preprocessing steps was described in the subsection “Experimental testing dataset” under the section “Performance comparisons among the models.” However, for greater specificity and reproducibility, in this Appendix, we also include the exact *afni_proc.py* command in AFNI (version AFNI_16.3.06) that was implemented to create the full processing pipeline. While there are several processing steps specified, each with many user-chosen options, it is possible to provide the exact pipeline in a succinct manner because the

processing steps and options were created and specified using *afni-proc.py* in AFNI. This tool permits the user full freedom to tailor a desired pipeline that may be reliably duplicated for the entire group, stored for future reference and published with a study for unambiguous description.

APPENDIX B

Data at three voxels from 25 subjects with two sessions that are illustrated¹² in Figure 4 and Table 4

Voxel Session Data	V_1				V_2				V_3			
	y_{1j}	$\hat{\sigma}_{1j}^2$	y_{2j}	$\hat{\sigma}_{2j}^2$	y_{1j}	$\hat{\sigma}_{1j}^2$	y_{2j}	$\hat{\sigma}_{2j}^2$	y_{1j}	$\hat{\sigma}_{1j}^2$	y_{2j}	$\hat{\sigma}_{2j}^2$
S1	0.075	0.013	-0.067	0.015	1.164	0.147	0.529	0.098	0.621	0.030	0.217	0.036
S2	0.160	0.006	0.081	0.004	0.705	0.016	0.100	0.026	0.407	0.018	-0.111	0.012
S3	0.101	0.006	0.084	0.009	-0.862	0.262	1.121	0.354	0.473	0.018	0.541	0.014
S4	-0.055	0.006	-0.041	0.004	0.297	0.050	0.113	0.014	0.012	0.008	-0.042	0.004
S5	0.075	0.010	0.053	0.010	0.416	0.031	0.412	0.044	0.036	0.011	-0.074	0.017
S6	0.282	0.009	0.199	0.010	0.590	0.012	0.166	0.019	0.906	0.011	0.539	0.014
S7	-0.012	0.014	0.273	0.010	1.042	0.038	0.727	0.030	0.486	0.012	0.506	0.008
S8	0.123	0.002	0.334	0.003	0.321	0.014	0.563	0.098	0.715	0.011	1.056	0.019
S9	-0.079	0.013	-0.213	0.013	0.366	0.012	-0.023	0.014	0.276	0.008	0.067	0.007
S10	-0.040	0.010	-0.127	0.015	0.705	0.026	0.734	0.028	0.439	0.015	0.331	0.020
S11	0.157	0.010	0.135	0.005	0.758	0.141	-0.033	0.159	0.280	0.013	-0.072	0.008
S12	0.056	0.004	-0.120	0.002	0.588	0.028	0.683	0.072	0.449	0.009	-0.052	0.005
S13	0.224	0.010	0.114	0.007	0.403	0.031	0.185	0.016	0.201	0.016	-0.294	0.015
S14	0.024	0.003	-0.007	0.003	0.713	0.016	0.501	0.021	0.453	0.011	0.061	0.013
S15	0.107	0.009	0.077	0.005	0.453	0.013	0.283	0.008	0.066	0.018	0.001	0.007
S16	0.182	0.008	0.160	0.006	1.120	0.020	0.303	0.017	0.438	0.011	0.019	0.006
S17	0.035	0.003	0.000	0.003	0.648	0.008	0.577	0.006	0.569	0.007	0.382	0.008
S18	0.014	0.003	-0.057	0.007	0.872	0.024	0.189	0.036	0.264	0.005	-0.025	0.009
S19	0.151	0.004	0.200	0.003	0.436	0.006	0.616	0.005	0.343	0.005	0.306	0.003
S20	0.165	0.004	0.020	0.003	0.278	0.007	0.386	0.009	0.398	0.004	0.434	0.007
S21	0.136	0.010	0.151	0.009	0.318	0.034	0.318	0.025	0.141	0.017	0.267	0.015
S22	0.165	0.012	0.177	0.021	-0.004	0.027	0.444	0.046	0.043	0.012	0.242	0.013
S23	0.005	0.005	0.085	0.003	0.410	0.005	0.217	0.016	0.322	0.004	0.313	0.005
S24	0.213	0.009	0.061	0.009	0.663	0.016	0.493	0.017	0.938	0.015	0.517	0.014
S25	0.059	0.010	0.132	0.008	1.191	0.011	0.318	0.012	0.324	0.008	0.001	0.011

¹²To save space, the data shown here are rounded to the nearest thousandth, therefore there may have some small differences with the results listed in Table 4.